# Search Engine for Software Packages

## Final Presentation - May 2nd, 2018

Joshua Choo, Mengshi Feng, Vivian Liu, Avery Nisbet, Yidan Zhang

# Agenda

- ❖ Project Overview
  - ➢ Industry Analysis
  - ➢ User Interviews
  - ➢ Current Solution vs Proposed Solution
- ❖ Project Components and Integration
- ❖ Demo
- ❖ Next Steps

# Search Engine Industry: Market Share (US)

**Major players**
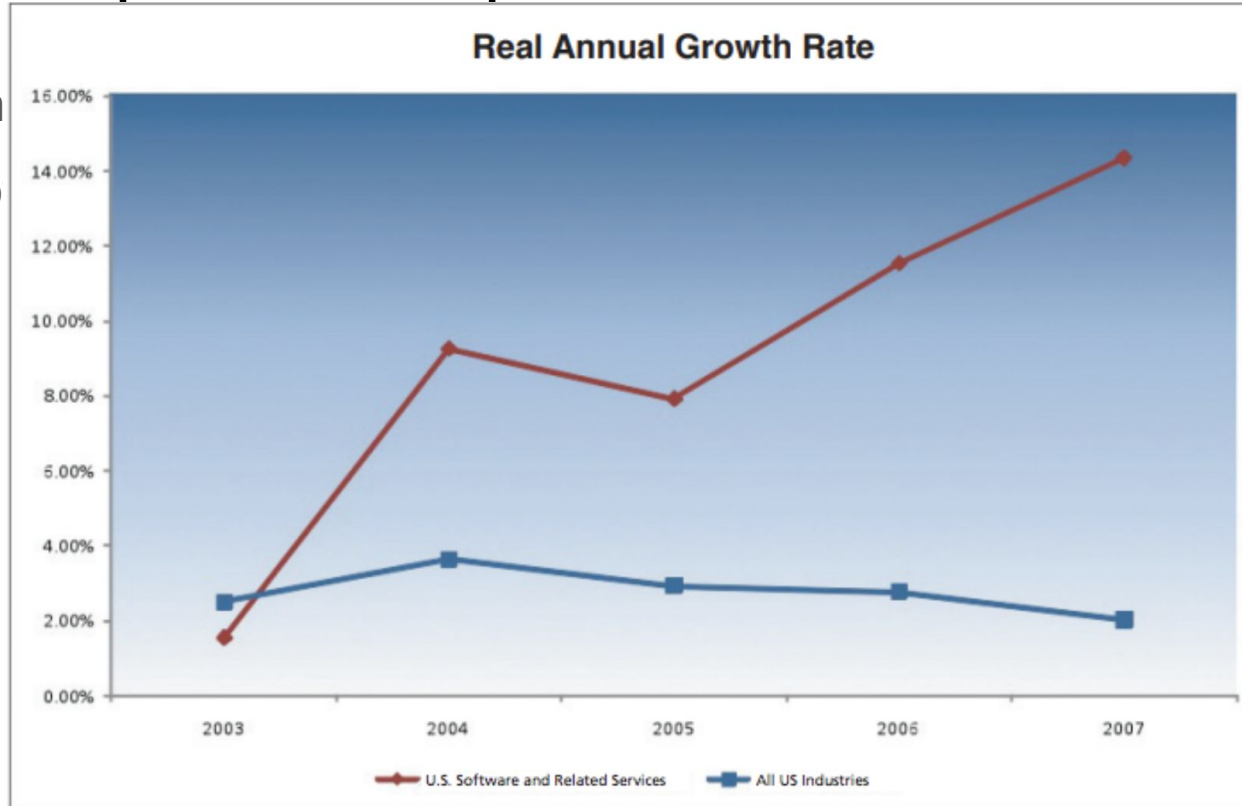(Market share)

Microsoft Corporation 11.1%

Alphabet Inc. 69.9%

19.0%
Other

SOURCE: WWW.IBISWORLD.COM

Google has a lion's share of the market

Search Engines in the US, IBISWorld Industry Report 51913a. (2017). IBISWorld, pp.4-34. Retrieved November 2, 2017, from
http://clients1.ibisworld.com/reports/us/industry/default.aspx?entid=1982

# Why is this problem important?

- 1 million
- 12 out o                                              stration

**Real Annual Growth Rate**



Legend: U.S. Software and Related Services — All US Industries

4

Business Software Alliance. (n.d.). Software Industry Facts and Figures. Retrieved November 4, 2017, from http://www.bsa.org/country/public%20policy/~/media/files/policy/security/general/sw_factsfigures.ashx
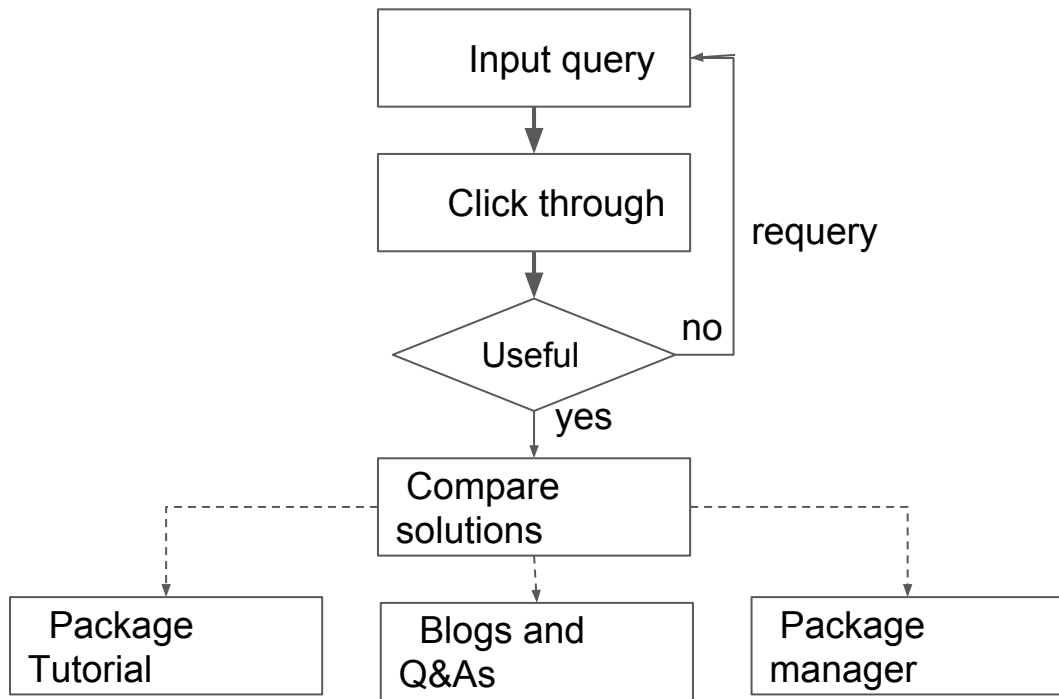
# Why is Software Search Difficult?

- Difficult to formulate precise queries
- Misleading or outdated documentation
- Time consuming to iteratively evaluate and aggregate useful information across links

# Software Package Search: A Problem Left to be Solved

Not a universal solution!

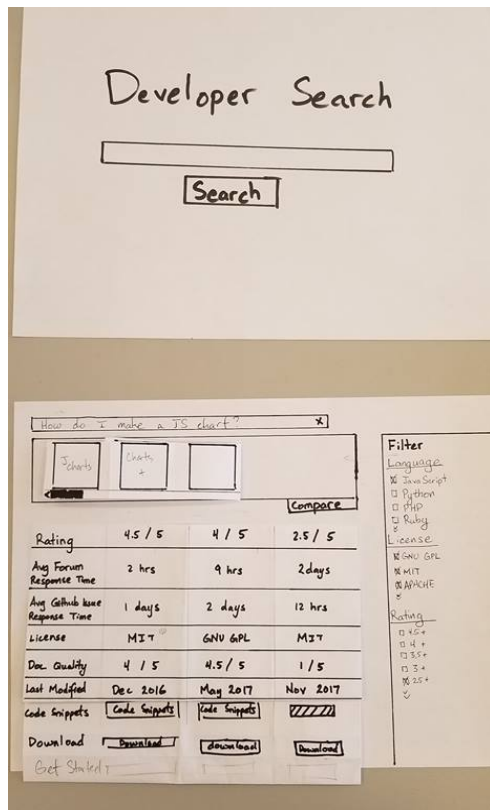# Setting Course: Narrowing Down the Solutions with User Interviews



Package evaluation depends on who searched, and what the search was



List of evaluation criteria created from interviews

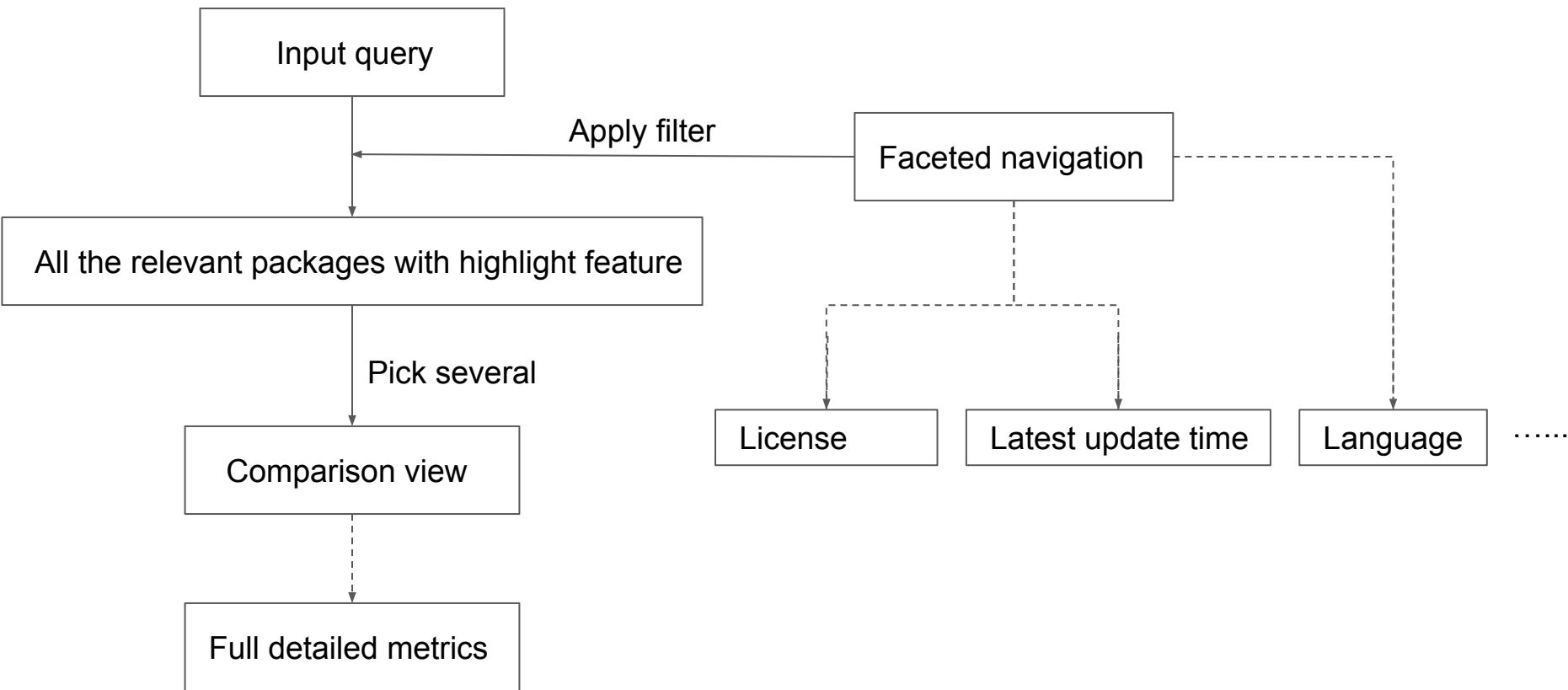(Gallardo-Valencia and Sim, 2011)

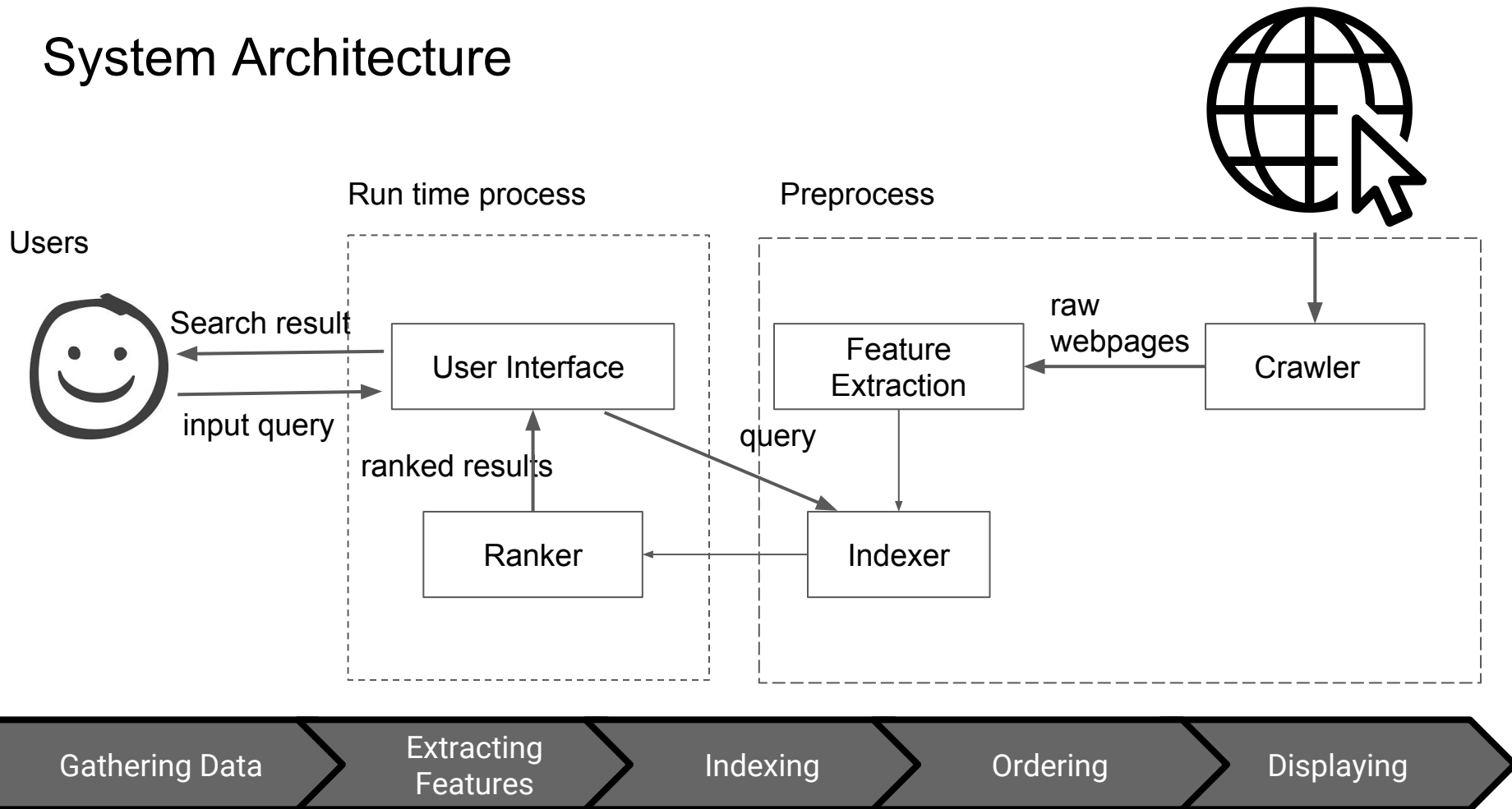# Setting Course: Narrowing Down the Solutions with User Interviews



Testing a low-fidelity prototype:

- Paper Prototypes made from Search Engine Design Features
- Users run through search experience
- Could easily draw a new feature to test

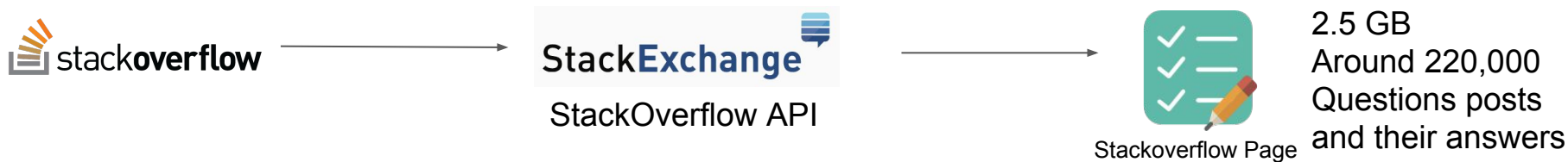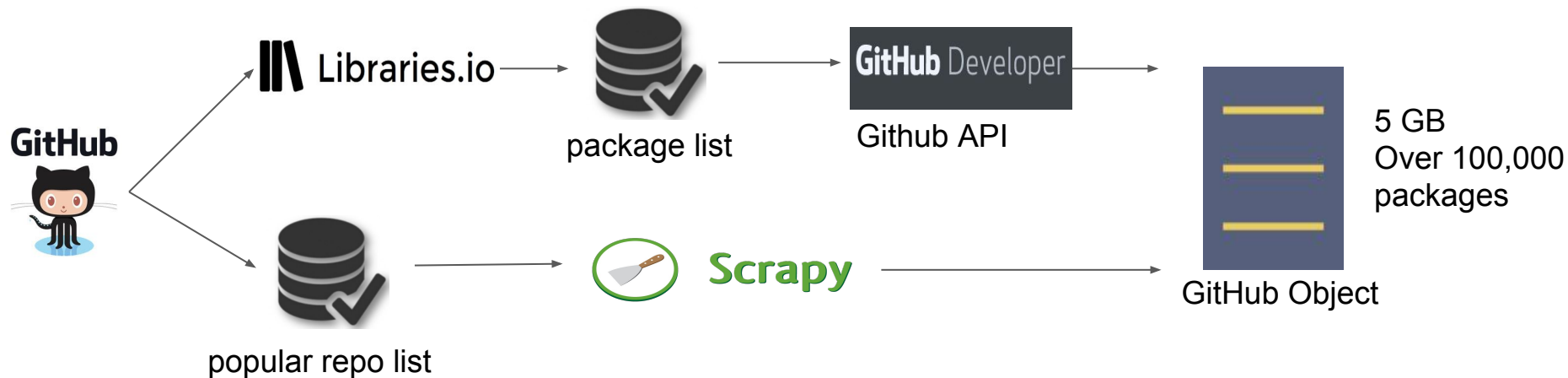# How will people search package using our system?

Input query

All the relevant packages with highlight feature

Comparison view

Full detailed metrics

Apply filter

Faceted navigation

Pick several

License

Latest update time

Language

······

# System Architecture



Run time process

Preprocess

Users

Search result

User Interface

input query

ranked results

Ranker

query

Indexer

Feature Extraction

raw webpages

Crawler

Gathering Data | Extracting Features | Indexing | Ordering | Displaying

# Crawler: Gathering Information



Libraries.io → package list → GitHub Developer / Github API → 5 GB Over 100,000 packages / GitHub Object

GitHub → popular repo list → Scrapy → GitHub Object

stackoverflow → StackExchange / StackOverflow API → Stackoverflow Page / 2.5 GB Around 220,000 Questions posts and their answers

| Gathering Data | Extracting Features | Indexing | Ordering | Displaying |

# Feature Extraction: Interpreting Unstructured Data



Python Scripts to convert crawled data

```
{
    "Name": "angular",
    "_childDocuments_": [
        {
            "Body_markdown": ".....",
            "Title": ".....",
            "Link": ".....",
            "Code_snippets": ".....",
            "Answers": ".....",
        },
        {
            "Body_markdown": ".....",
            "Title": ".....",
            "Link": ".....",
            "Code_snippets": ".....",
            "Answers": ".....",
        },
        ...........
    ]
}
```

GitHub Object — Stackoverflow Page

Gathering Data ▸ **Extracting Features** ▸ Indexing ▸ Ordering ▸ Displaying

# Indexer: Providing Quick Look-up on Data



Feature
Extraction

REST-ful api

Solr

django

Google Cloud

| Gathering Data | Extracting Features | Indexing | Ordering | Displaying |

# Ranker: Reordering indexed results

- Score relevance of packages based on Okapi BM25 similarities of package metadata (Github metadata, Stackoverflow documents and Slant queries) and query phrase.
- Each package is tagged with relevant Slant queries and similarity between them and the user's query are considered.
- Implemented slop parameters to ensure that relevant documents without the exact phrase query would be considered, while not inflating the scores of irrelevant documents that contains the individual words in the phrase query.

| Gathering Data | Extracting Features | Indexing | Ordering | Displaying |

# User Interface: Displaying Information Intuitively



Features:

- **Carousel** for packages.
- **Details** for the package appear on package selection.
- **Filters** help you narrow down your search.
- **Keywords** from the query are accepted as filters.
- **Compare** packages side by side.
- **Stackoverflow** results help you get opinions of other developers.
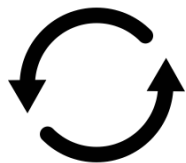
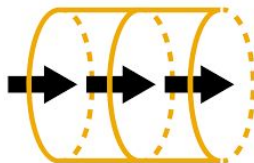| Gathering Data | Extracting Features | Indexing | Ordering | Displaying |

# Demo

http://35.230.66.167/

# Next Steps

**Continuous and Exhaustive Crawling**

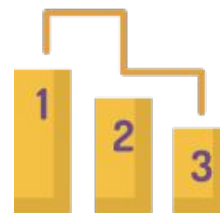**Data Pipeline**

**Recognition of User Queries**

**Extraction of More Metadata**

**Improved Ranking Algorithm**

# Question?



Users

Run time process

Preprocess

Search result

input query

User Interface

ranked results

Ranker

Indexer

query

Feature Extraction

raw webpages

Crawler

# Appendix -- Solr Schema

{

*1st level*

    "name": "Angular",

    "path": "1.git",

    "homepage_url": "https://angular.io",

    "headme": "...",

    ...

*2nd level*

    "_childDocuments_": {

        "Title": "Capture Video of Android&#39;s Screen",

        "path": "2.stack",

*3rd level*

        "_childDocuments_"

            "path": "3.stack.answer",

            "Answer_id": "x"

    }

}