# Search Engine

nutch

```
[Crawl] → [Crawling]
                 |
            [Spider] → [robot.txt]
```

[GFS] → [HDFS] → [BigData]

[Storage] → HBASE

[Pre processing] →
- → Parsing
- → lower case
- → punctuation
- → Special character
- → Stop words
- → Spell correction ⇐
- → Convert emoji
- → Root word ⇐

The
A
US ⇐

the
a
us

. , : ? — ~

Porter Stemming
rules

Stemming
grammatical rules
Lemmatization
↓
accurate
↓
slow

[Query] (red)
↓
[Correct] (red)
↓
[Retrieval] (red)

[Indexing]
↳ RDBS / NOSQL

[Structed Representation]

→ text → website NOUN
↑
(run) running runs
Sit Snt Sitting
Snt Stemming Sit
Stemm

Page Ranking ← (1000)
[text → website] (red)
Larry Page (red)
Video → (red)
Images → (red)
Ads → (red)

I had enough breakfast ⇒ | enough breakfast
breakfast was awesome ⇒ | breakfast awesome
I had to take a trip ⇒ | take trip
I went for a running ⇒ | go run

he runs very fast ⇒ | run very fast

he runs fast, while he ran he | run fast | run collabs
collapsed

n grams ⇒ Uni gram, bi grams, tri gram

|  | Uni gram | bi gram |
|---|---|---|
| enough breakfast → | [enough, breakfast] | [enough breakfast] |
| breakfast awesome → | [breakfast, awesome] | [breakfast aw] |
| take trip → | [take, trip] | [take trip] |
| go run ⇒ | [go, run] | [go run] |
| run very fast → | [run fast] | [run very] [very fast] |

| | enough | breakfast | awesome | take | trip | run enough break | break fast |
|---|---|---|---|---|---|---|---|
| 1. | 1 | 1 | 0 | 0 | 0 | 1 | 0 |
| 2. | 0 | 1 | 0 | 0 | 0 | 0 | 1 |
| 3. | | | | | | | |
| 4. | | | | | | | |
| 5. | | | | | | | |

# Indexing

1. One hot encoding →
   a. Occurances not there for documents and compus
   b. orders are not preserved
   c. Context is not present

2. Bag of Word
   a. Occurrences not there for corpus
   b. orders are not preserved
   c. Context is not there

3. TF-IDF

4. Word embedding