# Information Retrieval

## Venkatesh Vinayakarao

Term: Aug – Sep, 2019
Chennai Mathematical Institute

**Mission defines strategy, and strategy defines structure.**
**– Peter Drucker.**

# Query Understanding
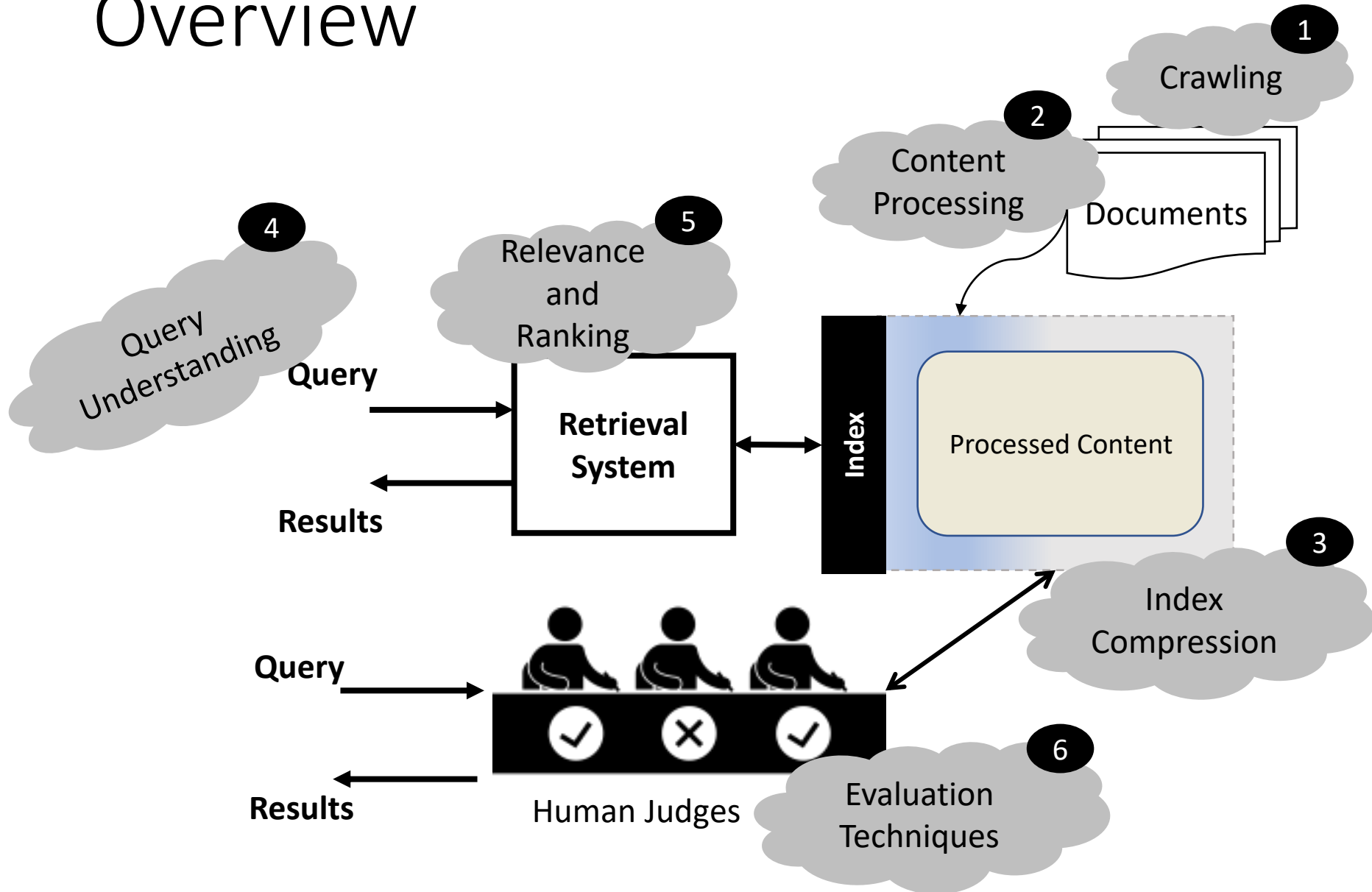
# Agenda

| An Overview of Query Types | Understanding the query types helps us to optimize the retrieval system. |
|---|---|
| Methods of Query Understanding | Token-level Query Processing (Query Segmentation, Spelling Correction, Phonetic Correction) |

# Overview

# Some Queries are Hard to Understand!

• Guess, what should the query "IR" return?

Google | IR

Q All    🗏 News    ⊙ Maps    🖾 Images   ⊘

About 3,03,00,00,

www.indianrail.gov
**Indian Railwa**

irctc.co.in ▾
**Irctc**

b | IR

ALL    IMAGES    VIDEOS    MAPS    NEWS

17,70,00,000 Results    Date ▾    Language ▾    Region ▾

**Welcome to Indian Railway Passenger Reservation Enquiry**
www.indianrail.gov.in ▾
The official site with information on trains, fares and availability.
PNR Status · National Train Enquiry System

**Videos of ir**
bing.com/videos

Depends on the context: who is querying, when they are querying and what was queried before, popularity of keyword, etc.

# Query Types: The N, I & T!

- Navigational
  - Example: "fb"
  - Say, a user wants to visit facebook.com. He might hit fb on the search bar and use the first result to go to the page.

- Informational
  - Example: "Amitabh bachchan"
  - Seeks information about Amitabh Bachchan.

- Transactional
  - Example: "Chennai to Delhi air ticket"
  - Say, the intent is to buy an air ticket and this query is the first step of searching for the best price/route/vendor.

# Query Types: Long and Short

- Typically, queries are **short phrases**
  - "data science degree india"
  - "Stanford semester start date"
- However, **long queries** are not uncommon
  - Example: "easter egg hunts in northeast columbus parks and recreation centers"
  - "Queries of length five words or more have increased at a year over year rate of 10%, while single word queries dropped 3%." – Balasubramanian, Kumaran and Carvalho – 2010.

# Query Types: Head and Tail Queries

- Head Queries
  - Queries that appear very frequently

- Tail Queries
  - The "rare" queries

In a quest to improve overall performance, we often do not give the attention here that this deserves

# Query Types: Question and Answers



Wow! How did Bing understand this?

Good Query Understanding

# Methods for Query Understanding

- Token-level Query Processing
  - Spelling Errors
  - Query Segmentation
- Query Reduction
  - Remove less-important query tokens.
- Query Expansion
  - Add more terms to query to improve precision and recall.
- Query Rewriting
  - Transform the original query to a query friendlier to the retrieval system.

# Token-Level Query Processing

# Query Segmentation

- Users might miss spaces when they query. Consider, for example:
  - Statebankofindia for "State Bank of India"
  - Amazonprimevideo for "Amazon Prime Video"
- Can you give an algorithm to check if input can be split to arrive at dictionary terms?

# A Recursive Algorithm

| S | T | A | T | E | B | A | N | K | O | F | I | N | D | I | A |

| S | █ | T | A | T | E | B | A | N | K | O | F | I | N | D | I | A |

S not in our dictionary.  Keep moving till we find a dictionary term.

| S | T | A | T | E | █ | B | A | N | K | O | F | I | N | D | I | A |

Check if rest of the string "BANKOFINDIA" can be split.
If "yes", insert a space after STATE.
Else, continue with a longer term.

| S | T | A | T | E | █ | B | A | N | K | █ | O | F | █ | I | N | D | I | A |

Recurse and backtrack till you find a split.

| **Dictionary** |
| --- |
| State |
| Bank |
| Of |
| India |
| Amazon |
| Prime |
| Video |
| … |

# A Python Solution

```python
def splitQuery(dictionary, str):
    if not str:
        return True

    for i in range(1, len(str) + 1):
        prefix = str[:i]
        if prefix in dictionary and splitQuery(dictionary, str[i:]):
            return True
    return False

if __name__ == '__main__':
    dictionary = ["state", "bank", "india", "of", "amazon", "prime", "video"]
    query = "statebankofindians"
    if splitQuery(dictionary, query):
        print("Yes, Query can be split into terms.")
    else:
        print("No. Query cannot be split.")
```

# Spelling Errors

# Some Types of Misspellings

| Cause | Misspelling | Correction |
|---|---|---|
| Typing quickly | exxit mispell | exit misspell |
| Keyboard adjacency | importamt | important |
| Inconsistent rules | concieve conceirge | conceive concierge |
| Ambiguous word breaking | silver light | silverlight |
| New words | kinnect | kinect |

# Spelling Errors in Query

- 10% to 20% of queries carry misspelt words[1].

- English is not 100% phonetic (e.g., colonel, read vs dead).

- How many of these phrases contain spelling errors?
    - cigarete lighter
    - fourty dollars
    - going to libary today
    - unforgetable holiday
    - successful businessman

Duan and Hsu, WWW 2011.

# Notorious Britney

The data below shows some of the misspellings detected by our spelling correction system for the query [ britney spears ], and the count of how many different users spelled her name that way.  -- Google.

```
488941 britney spears      29 britent spears      9 brinttany spears     5 brney spears       3 britiy spears        2 brirreny spears
 40134 brittany spears     29 brittnany spears     9 britanay spears      5 broitney spears    3 britmeny spears      2 brirtany spears
 36315 brittney spears     29 britttany spears     9 britinany spears     5 brotny spears      3 britneeey spears     2 brirttany spears
 24342 britany spears      29 btiney spears        9 britn spears         5 bruteny spears     3 britnehy spears      2 brirttney spears
  7331 britny spears       26 birttney spears      9 britnew spears       5 btiyney spears     3 britnely spears      2 britain spears
  6633 briteny spears      26 breitney spears      9 britneyn spears      5 btrittney spears   3 britnesy spears      2 britane spears
  2696 britteny spears     26 brinity spears       9 britrney spears      5 gritney spears     3 britnetty spears     2 britaneny spears
  1807 briney spears       26 britenay spears      9 brtiny spears        5 spritney spears    3 britnex spears       2 britania spears
  1635 brittny spears      26 britneyt spears      9 brtittney spears     4 bittny spears      3 britneyxxx spears    2 britann spears
  1479 brintey spears      26 brittan spears       9 brtny spears         4 bnritney spears    3 britnity spears      2 britanna spears
  1479 britanny spears     26 brittne spears       9 brytny spears        4 brandy spears      3 britntey spears      2 britannie spears
  1338 britiny spears      26 btittany spears      9 rbitney spears       4 brbritney spears   3 britnyey spears      2 britannt spears
  1211 britnet spears      24 beitney spears       8 birtiny spears       4 breatiny spears    3 britterny spears     2 britannu spears
  1096 britiney spears     24 birteny spears       8 bithney spears       4 breetney spears    3 brittneey spears     2 britanyl spears
   991 britaney spears     24 brightney spears     8 brattany spears      4 bretiney spears    3 brittnney spears     2 britanyt spears
   991 britnay spears      24 brintiny spears      8 breitny spears       4 brfitney spears    3 brittnyey spears     2 briteeny spears
   811 brithney spears     24 britanty spears      8 breteny spears       4 briattany spears   3 brityen spears       2 britenany spears
   811 brtiney spears      24 britenny spears      8 brightny spears      4 brieteny spears    3 briytney spears      2 britenet spears
   664 birtney spears      24 britini spears       8 brintay spears       4 briety spears      3 brltney spears       2 briteniy spears
   664 brintney spears     24 britnwy spears       8 brinttey spears      4 briitny spears     3 broteny spears       2 britenys spears
   664 briteney spears     24 brittni spears       8 briotney spears      4 briittany spears   3 brtaney spears       2 britianey spears
   601 bitney spears       24 brittnie spears      8 britanys spears      4 brinie spears      3 brtiiany spears      2 britin spears
   601 brinty spears       21 biritney spears      8 britley spears       4 brinteney spears   3 brtinay spears       2 britinary spears
   544 brittaney spears    21 birtany spears       8 britneyb spears      4 brintne spears     3 brtinney spears      2 britmy spears
   544 brittnay spears     21 biteny spears        8 britnrey spears      4 britaby spears     3 brtitany spears      2 britnaney spears
```

Source: https://archive.google.com/jobs/britney.html

# Two Major Approaches

- Two major approaches exist for spelling correction:
  - finding "**nearest**" dictionary term.
  - finding "**most commonly used**" dictionary term when there are multiple "nearest terms".
- Two major kinds:
  - Isolated-Term Correction
    - Correct one word at a time.
  - Context-Sensitive Correction
    - "flew *form* New York" – Note that form is a dictionary term. Yet, this requires to be corrected to "flew from New York".

# Edit Distance

# Spelling Correction: Edit distance

- Given two strings $S_1$ and $S_2$, the minimum number of operations to convert one to the other

- Operations are typically character-level
  - Insert, Delete, Replace, (and perhaps Transposition*)

- E.g., the edit distance from **dof** to **dog** is 1
  - From **cat** to **act** is 2          (Just 1 with transpose.)
  - from **cat** to **dog** is 3.

*In this course, we do not consider transposition.

# Quiz

What is the edit distance between Sunday and Saturday?

*You are allowed to perform only Insert, Delete, and Replace operations.

# Answer

- Saturday = Sunday = S*day

- Problem is same as
  - What is the edit distance between atur and un?
  - Answer
    - Delete a,t. Replace r with n.
    - 3.

# Levenshtein Example

|   |   | s | a | t | u | r | d | a | y |
|---|---|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| s | 1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| u | 2 | 1 | 1 | 2 | 2 | 3 | 4 | 5 | 6 |
| n | 3 | 2 | 2 | 2 | 3 | 3 | 4 | 5 | 6 |
| d | 4 | 3 | 3 | 3 | 3 | 4 | 3 | 4 | 5 |
| a | 5 | 4 | 3 | 4 | 4 | 4 | 4 | 3 | 4 |
| y | 6 | 5 | 4 | 4 | 5 | 5 | 5 | 4 | 3 |

**Sunday**

**Keep s. Insert a, t.**

**Keep u.**

**Replace r.**

**Keep day.**

**Saturday**

$$
\begin{aligned}
D(i,j) &= \min[D(i-1,j) + w_d, \\
&\quad\quad\; D(i,j-1) + w_i, \\
&\quad\quad\; D(i-1,j-1) + w_r] \quad \left.\right\} \forall i,j > 0 \\
D(i,0) &= D(i-1,0) + w_d \\
D(0,j) &= D(0,j-1) + w_i
\end{aligned}
$$

$$
D(0,0) = 0
$$

V. I. Levenshtein, Binary codes capable of correcting deletions insertions and reversals. Soviet Physics. 10, 707-710, 1966.
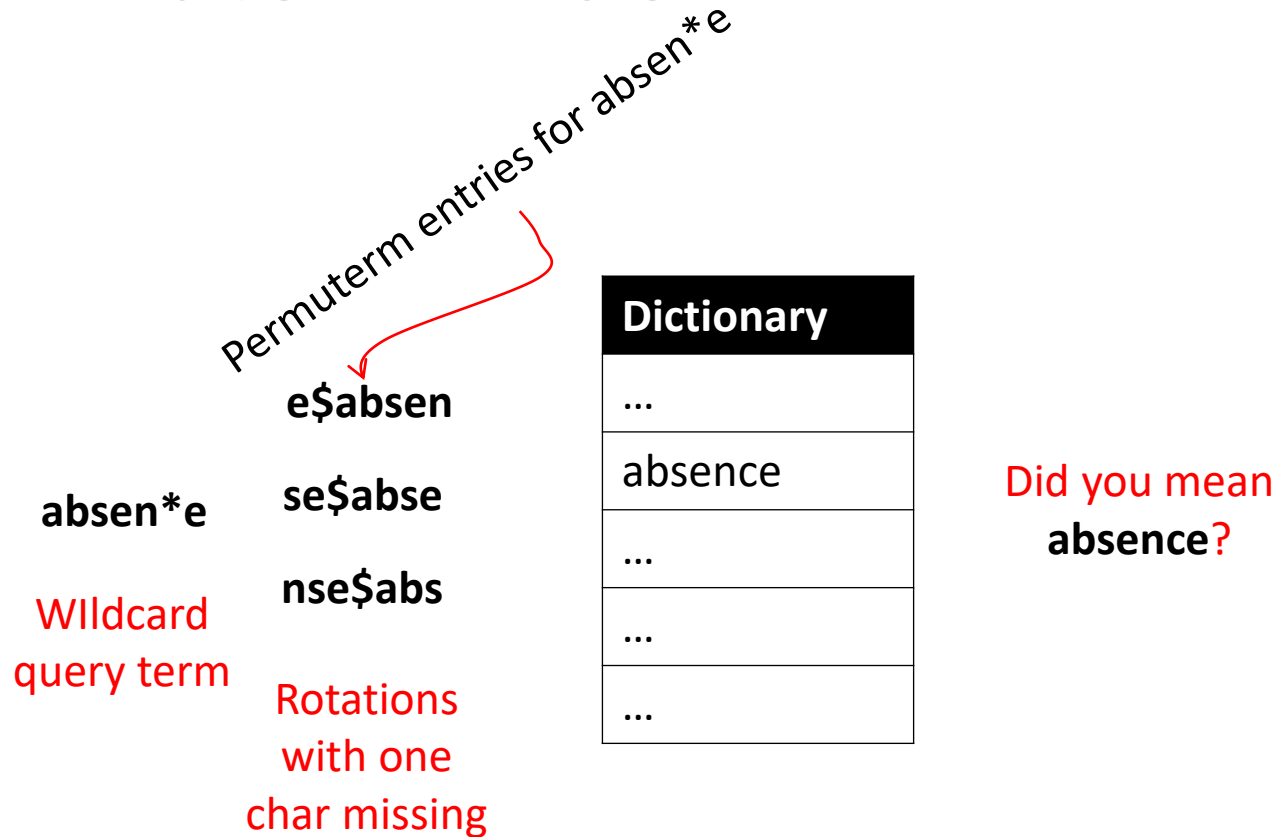
# Levenshtein Algorithm

EDITDISTANCE$(s_1, s_2)$

1   $int \ m[i,j] = 0$
2   **for** $i \leftarrow 1$ **to** $|s_1|$
3   **do** $m[i,0] = i$
4   **for** $j \leftarrow 1$ **to** $|s_2|$
5   **do** $m[0,j] = j$
6   **for** $i \leftarrow 1$ **to** $|s_1|$
7   **do for** $j \leftarrow 1$ **to** $|s_2|$
8       **do** $m[i,j] = \min\{m[i-1,j-1] + \text{if } (s_1[i] = s_2[j]) \text{ then } 0 \text{ else } 1\text{fi},$
9                          $m[i-1,j] + 1,$
10                         $m[i,j-1] + 1\}$
11  **return** $m[|s_1|, |s_2|]$

**Can we use Levenshtein's Distance to answer wildcard queries?**

# Permuterm Index

Permuterm entries for absen*e

**e$absen**

**absen*e**

**se$abse**

WIldcard query term

**nse$abs**

Rotations with one char missing

| Dictionary |
| --- |
| ... |
| absence |
| ... |
| ... |
| ... |

Did you mean **absence**?

---

Compute edit distance for each query term with each of its permuterm based matches. Very expensive!
Assume first letter will be correct. Apply such heuristics.

# K-grams for Spelling Correction

# k-gram Idea for Spelling Correction

- Many heuristics lead to poor matches.

- For example, "bored" misspelt as "bord" may match "boardroom" if the heuristic is
  - Match any two bigrams
    - and we matched "bo" and "rd"

- Potential Solution
  - Compute Jaccard Similarity between k-grams of matched term and that of the query term.

# Jaccard Coefficient

- Jaccard Coefficient of two sets A and B
  
  = |A∩B|/|A∪B|

- Example: JS on bigrams of ("bord", "boardroom")
  
  = |{$b, bo, rd}|/|{$b, bo, or, rd, d$, oa, ar, dr, ro, oo, om, m$}|
  
  = 3/12.

*If you do not use end markings, we get 2/9.

# Context-Sensitive Spelling Correction

- Our heuristics may lead to
  - "flew form Delhi" → "flew fore Delhi", "flew from Delhi"
- Surrounding words may determine the correction
- Potential Solution
  - Use query log frequency or collection frequency of these phrases to choose the best.

# Thank You