

Information Retrieval Engines

Seung-Hoon Na

Chonbuk National University

2017-05-04

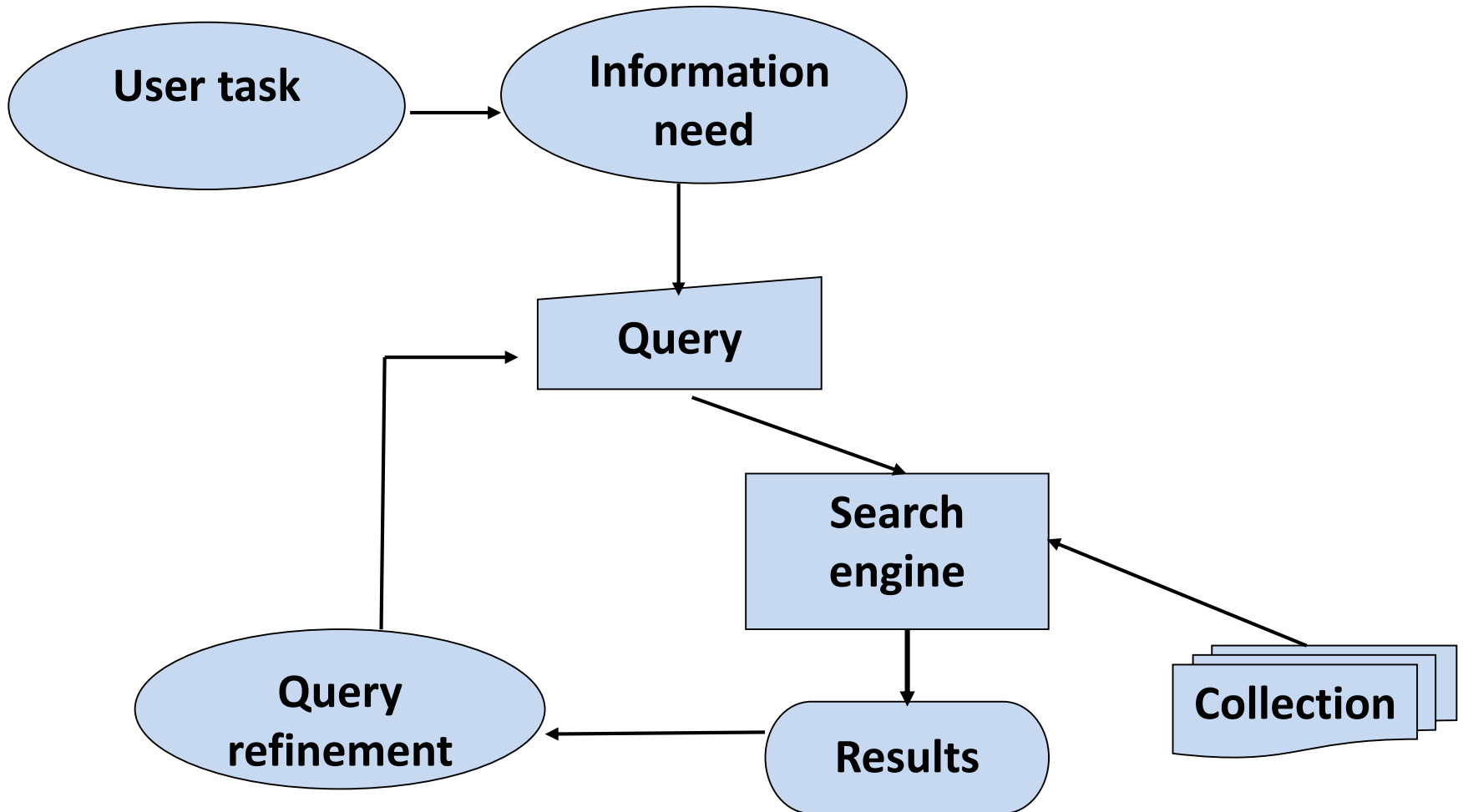
<https://nlp.stanford.edu/IR-book/>

Many slides from <http://nlp.stanford.edu/IR-book/newslides.html>

Information Retrieval: Introduction

- **Collection:** A set of documents
 - Assume it is a static collection for the moment
- **Goal:** Retrieve documents with information that is **relevant** to the user's **information need** and helps the user complete a **task**

Information Retrieval System



Boolean retrieval

- The Boolean model is the simplest model to base an information retrieval system on.
- Queries are Boolean expressions
 - 예) $Q = \{\text{전북대, 컴공}\}$
- Documents are considered as a **set** of terms.
- The search engine returns all documents that satisfy the Boolean expression.

Boolean retrieval:

Binary term-document incidence matrix

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	1	1	0	0	0	1
Brutus	1	1	0	1	0	0
Caesar	1	1	0	1	1	1
Calpurnia	0	1	0	0	0	0
Cleopatra	1	0	0	0	0	0
mercy	1	0	1	1	1	1
worser	1	0	1	1	1	0

Each document is represented by a binary vector $\in \{0,1\}^{|M|}$

Entry is 1 if term occurs, Example: CALPURNIA occurs in *Julius Caesar*.

Entry is 0 if term doesn't occur. Example: CALPURNIA doesn't occur in *The tempest*.

The vector space model

Queries & documents as Term weight vectors

	Anthony and Cleopatra	Julius C aesar	The Tempest	Hamlet	Othello	Macbeth ...
ANTHONY	5.25	3.18	0.0	0.0	0.0	0.35
BRUTUS	1.21	6.10	0.0	1.0	0.0	0.0
CAESAR	8.59	2.54	0.0	1.51	0.25	0.0
CALPURNIA	0.0	1.54	0.0	0.0	0.0	0.0
CLEOPATRA	2.85	0.0	0.0	0.0	0.0	0.0
MERCY	1.51	0.0	1.90	0.12	5.25	0.88
WORSER	1.37	0.0	0.11	4.15	0.25	1.95
...						

- Each document is now represented as a real-valued vector of **term weights** $\in \mathbb{R}^{|V|}$.

Cosine similarity between query and document

$$\cos(\vec{q}, \vec{d}) = \text{SIM}(\vec{q}, \vec{d}) = \frac{\vec{q} \cdot \vec{d}}{|\vec{q}| |\vec{d}|} = \frac{\sum_{i=1}^{|\mathcal{V}|} q_i d_i}{\sqrt{\sum_{i=1}^{|\mathcal{V}|} q_i^2} \sqrt{\sum_{i=1}^{|\mathcal{V}|} d_i^2}}$$

- q_i is the term weight of term i in the query.
- d_i is the term weight of term i in the document.
- $|\vec{q}|$ and $|\vec{d}|$ are the lengths of \vec{q} and \vec{d} .
- This is the **cosine similarity** of \vec{q} and \vec{d} or, equivalently, the cosine of the angle between \vec{q} and \vec{d} .

Documents as Term weight vectors

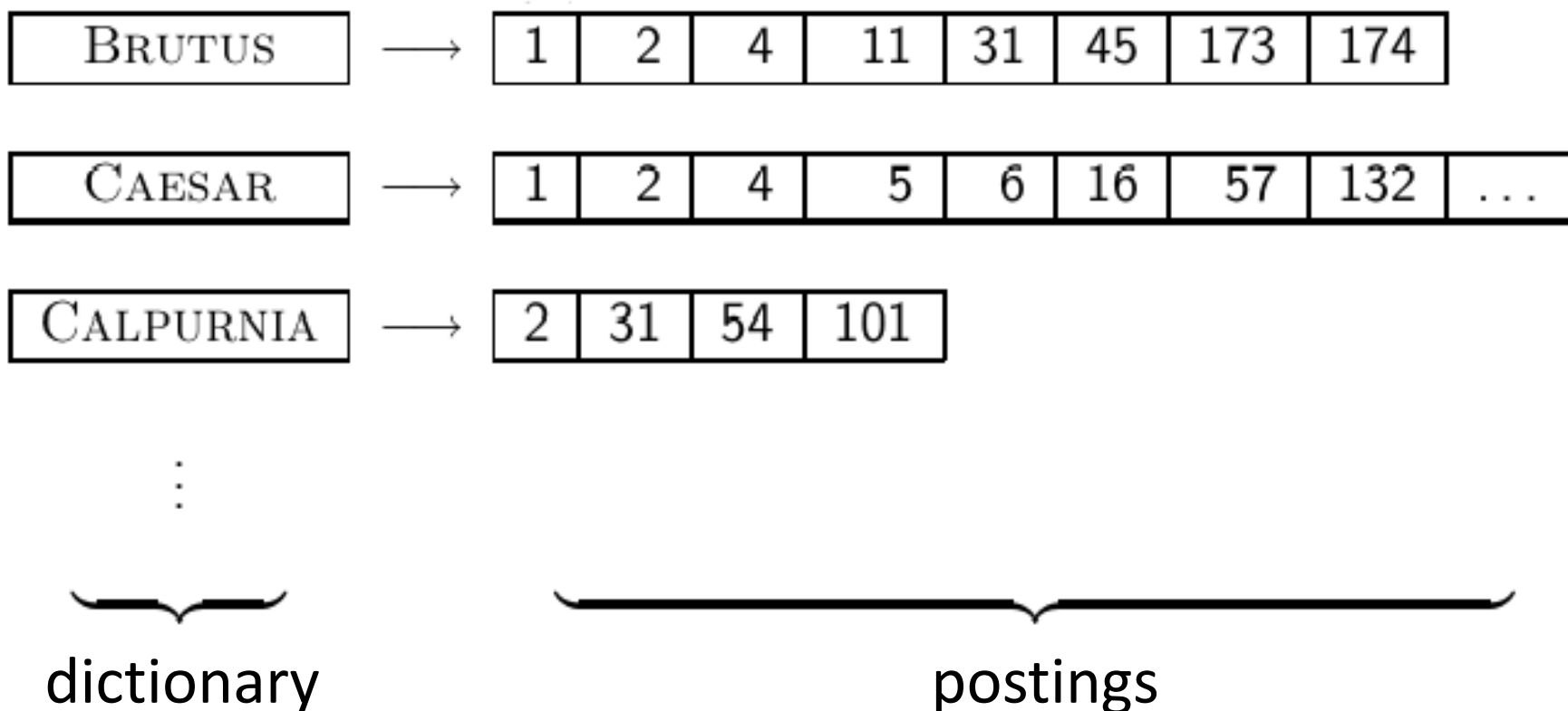
- $\vec{d} = \langle d_{w_1}, \dots, d_{w_t}, \dots, d_{|V|} \rangle$
 - d_{w_t} : document term weight for w_t
 - $d_{w_t} = c(w_t, d) * idf(w_t)$
 - $c(w_t, d)$: term frequency
 - 문서 (document d)내 용어 (term) 출현 빈도수
 - $idf(w_t)$: inverse document frequency
 - $idf(w_t) = \log \frac{N}{df(w_t)}$
 - $df(w_t)$: the document frequency, the number of documents that w_t occurs in
 - N : the number of total documents

Queries as Term weight vectors

- $\vec{q} = \langle q_{w_1}, \dots, q_{w_t}, \dots, q_{|V|} \rangle$
 - q_{w_t} : document term weight for w_t
 - $q_{w_t} = c(w_t, q)$
 - $c(w_t, q)$: query term frequency
 - 질의 (query)내 용어 (term) 출현 빈도수

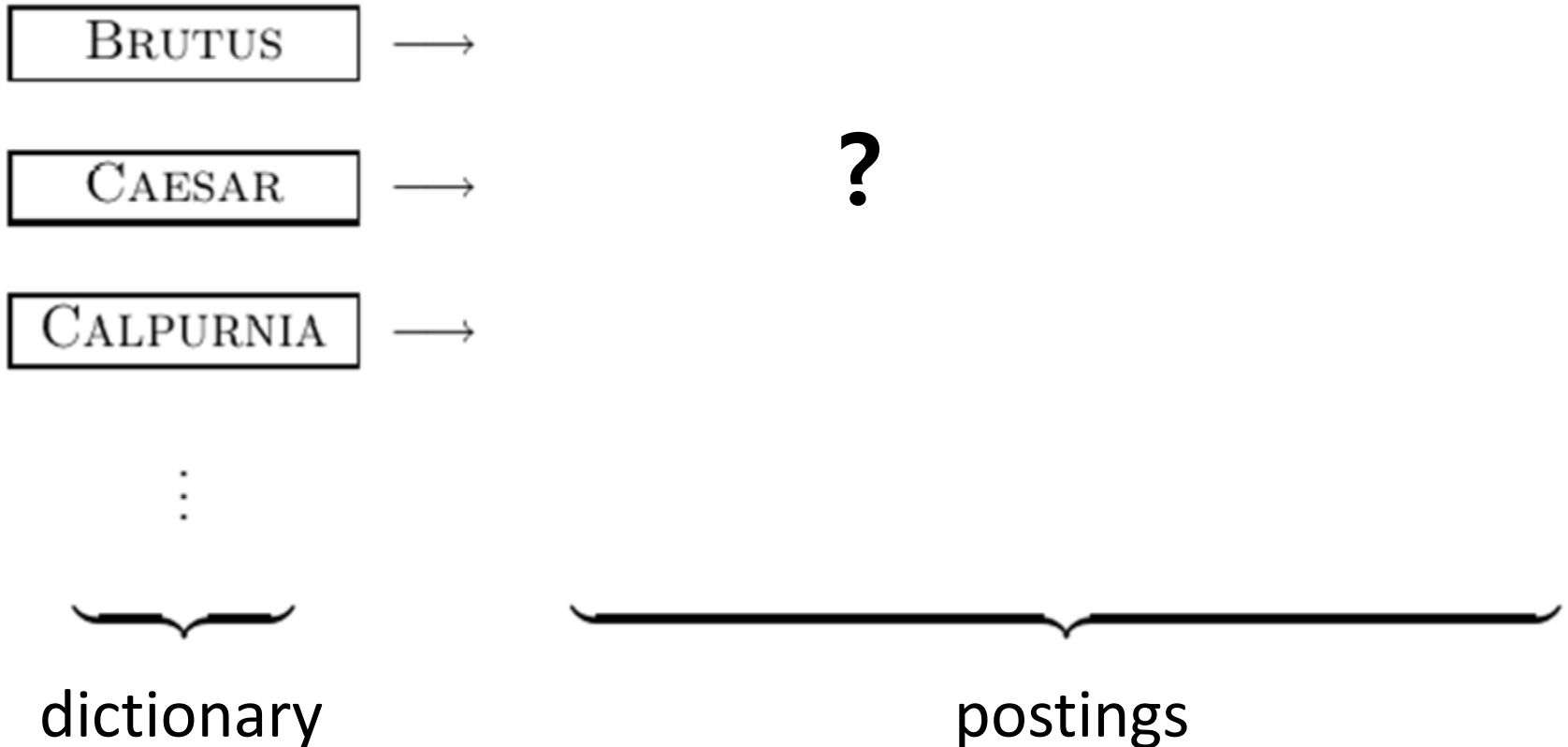
Inverted Index

For each term t , we store a list of all documents that contain t .



Inverted Index (with term frequency)

For each term t , we store a list of all documents that contain t .



Computing the Cosine score

COSINESCORE(q)

```
1  float Scores[ $N$ ] = 0
2  float Length[ $N$ ]
3  for each query term  $t$ 
4  do calculate  $w_{t,q}$  and fetch postings list for  $t$ 
5     for each pair( $d, tf_{t,d}$ ) in postings list
6     do Scores[ $d$ ] + =  $w_{t,d} \times w_{t,q}$ 
7  Read the array Length
8  for each  $d$ 
9  do Scores[ $d$ ] = Scores[ $d$ ]/Length[ $d$ ]
10 return Top  $K$  components of Scores[]
```

Term Extraction: character bi-gram

- 어절별로 char bigram 추출
 - 예) “전북대 컴퓨터공학부”



{전북, 북대, 컴퓨, 퓨터, 터공, 공학, 학부}

- 어절이 한 음절로 구성된 경우: unigram 사용
 - 예) “**몇** 시에” → {몇, 시에}

Assignment: 요약

- 1. 색인기 (indexer)
 - Inverted index 구조를 설계
 - 주어진 문서집합에 대한 inverted index 및 document content (텍스트 정보)를 db에 저장
 - Document content: Docid, Title, Text로 구성됨
- 2. 검색기 (searcher)
 - 주어진 Query에 대해서 상위 k개의 문서를 검색해주는 프로그램
 - 랭크, 문서번호, 문서제목(title), cosine score 를 순서대로 출력

Assignment: 검색기

- 주어진 질의(query)에 대해 상위 top K 문서 검색
- 출력 포맷:

총 num_rets 개의 검색 결과
(About num_rets results)

검색된 총 문서 결과 수 (num_rets) 출력 필요
1) 콘솔모드: 화면 하단에 출력
2) GUI모드: 화면 상단에 출력

- 1: $docname_1$ $score_1$ $title_1$ ← Rank-1 검색 문서 정보
- 2: $docname_2$ $score_2$ $title_2$ ← Rank-2 검색 문서 정보
- ...
- 50: $docname_{50}$ $score_{50}$ $title_{50}$

검색 결과가 많아 한 화면에 제시가 안될 경우 page view 지원

Assignment: 검색기

- 다음 두 가지 모드로 구현
 - 1) 콘솔모드
 - 2) swingx기반 GUI 모드
- 전체 검색 문서 수 (num_rets) 출력
 - 주어진 질의에 대해 전체 검색문서 수 num_rets도 함께 출력
- Page view지원
 - 전체 검색 문서 수 (num_ret)가 많아 결과가 한 화면에 출력
이 안될시에 page당 50개씩 출력하도록 page view 지원

Assignment: Dataset

- 문서포맷

```
<DOC>
<DOCNAME>2008102760218</DOCNAME>
<TITLE>안심했던 IC카드도 복제된다</TITLE>
<DATE>OCTOBER 27, 2008 09:10</DATE>
<TEXT>
복제 가능성이 제기돼 논란을 빚고 있는 집적회로(IC)칩 내장 현금카드가 정부 연구소의 실험에서 실제로 복제된 사실이 있었던 것으로 확인>됐다.
한나라당 진수희 의원이 26일 입수한 IC현금카드 복제결과 보고서에 따르면 한국전자통신연구원(ETRI) 산하 국가보안기술연구소(NSRI)는 8월 IC현금카드에 대한 복제 실험을 실시해 내장 암호키 추출 및 카드 복제에 성공한 것으로 나타났다. ...
</TEXT>
</DOC>
```

문서ID

제목

날짜

본문

Assignment

- Deadline: 6월 19일
- 데이터셋 Download: 메일 참조
- 제출 방법: 메일 참조
- 상세 추가 요구사항
 - 1) 색인기는 콘솔모드
 - 2) 검색기는 콘솔/GUI 모드: Page view 지원
 - 3) 결과 보고서 제출
 - 1. 정보검색 방법 개요
 - 2. 객체 설계도 및 구현 상세 내용
 - 3. 수행 결과 (색인/검색 결과 예)
 - 검색 결과는 콘솔뷰 및 gui뷰 각각 제시