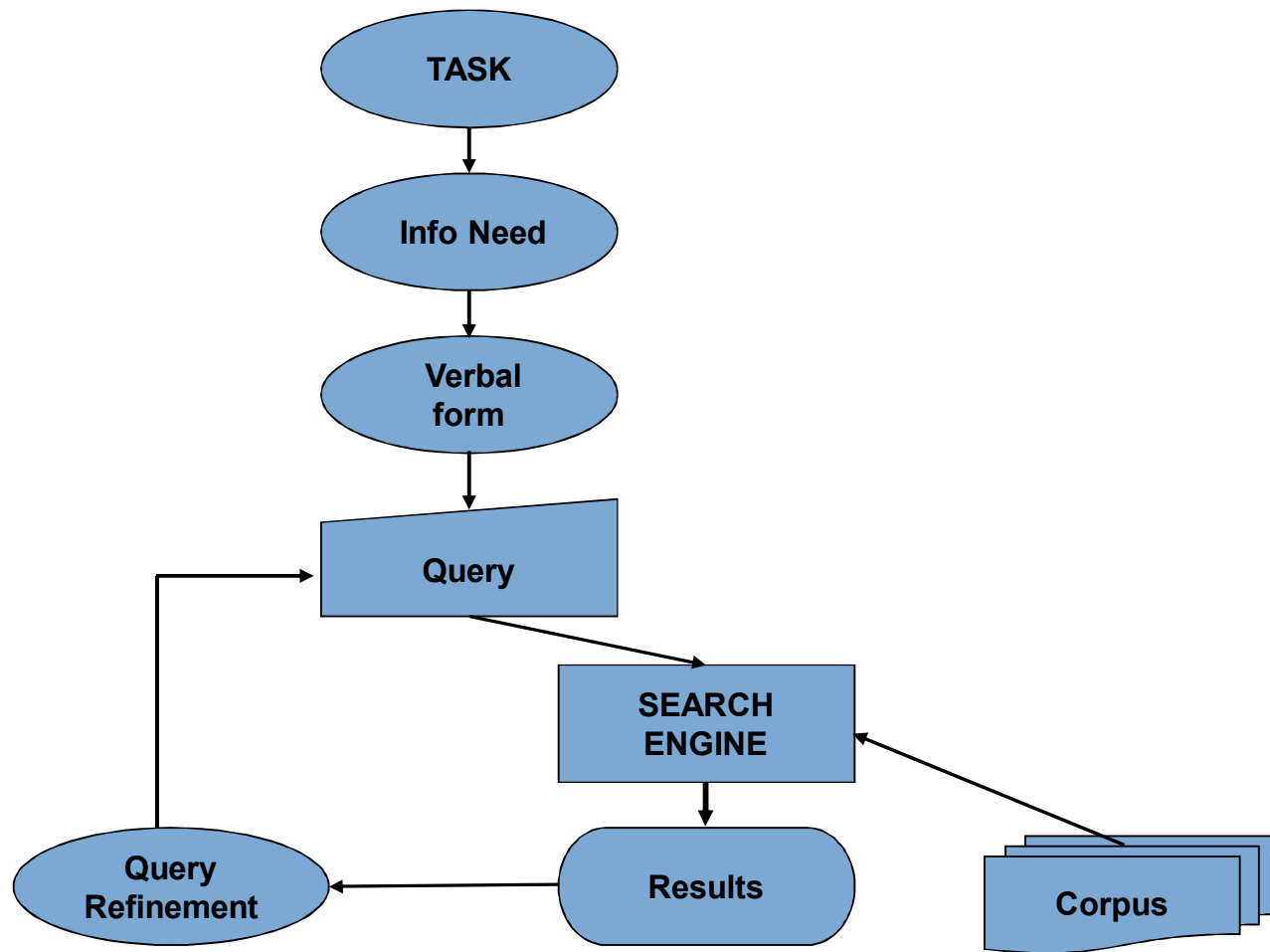


Search Engine (Basic concept)

The classic search model



Inverted index (역색인) construction

Documents to
be indexed



Friends, Romans, countrymen.

Tokenizer

Token stream

Friends

Romans

Countrymen

Linguistic modules

Modified tokens

friend

roman

countryman

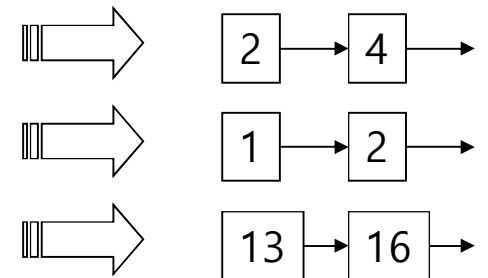
Indexer

Inverted index

friend

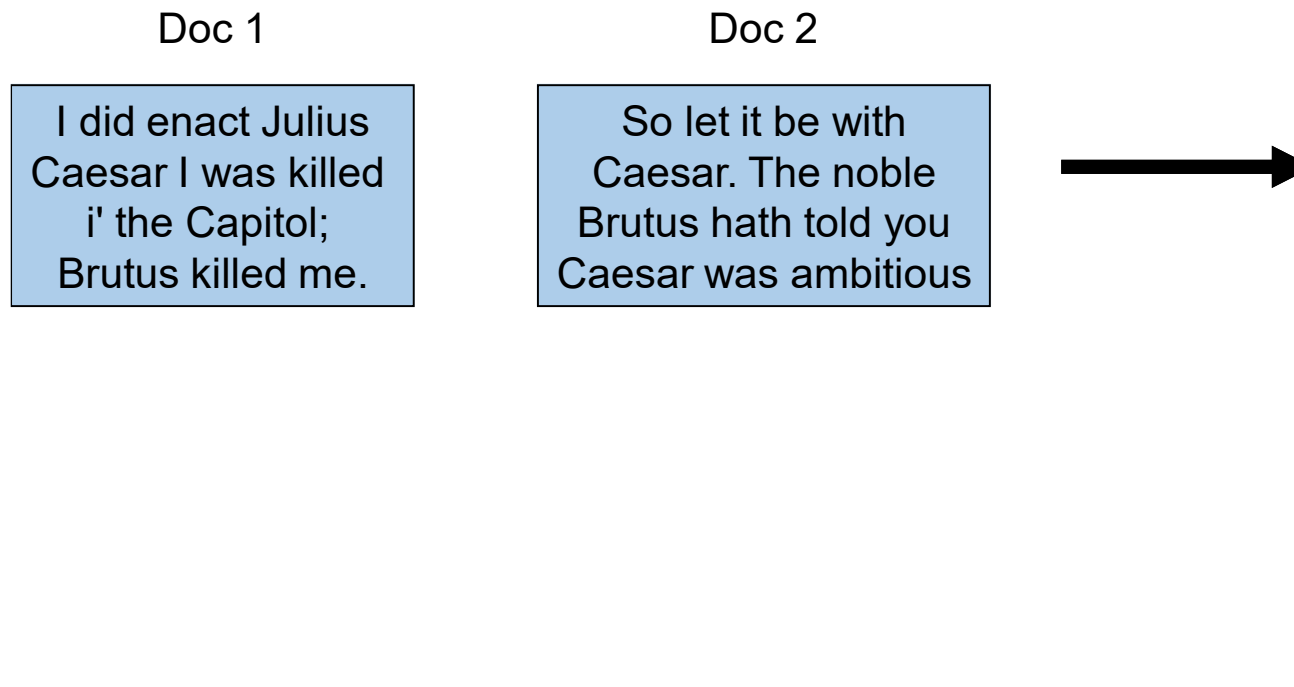
roman

countryman



Indexer steps: Token sequence

- Sequence of (Modified token, Document ID) pairs.



Indexer steps: Sort

- Sort by terms
 - And then docID

Core indexing step

Term	docID
I	1
did	1
enact	1
julius	1
caesar	1
I	1
was	1
killed	1
i'	1
the	1
capitol	1
brutus	1
killed	1
me	1
so	2
let	2
it	2
be	2
with	2
caesar	2
the	2
noble	2
brutus	2
hath	2
told	2
you	2
caesar	2
was	2
ambitious	2



Term	docID
ambitious	2
be	2
brutus	1
brutus	2
capitol	1
caesar	1
caesar	2
caesar	2
did	1
enact	1
hath	1
I	1
I	1
i'	1
it	2
julius	1
killed	1
killed	1
let	2
me	1
noble	2
so	2
the	1
the	2
told	2
you	2
was	1
was	2
with	2

Indexer steps: Dictionary & Postings

- 동일 문서에서 같은 용어가 여러 번 출현하면 하나로 합침
- 사전과 포스팅으로 분리
- 사전에 문서 빈도 (Doc. Frequency)를 저장

Term	docID
ambitious	2
be	2
brutus	1
brutus	2
capitol	1
caesar	1
caesar	2
caesar	2
did	1
enact	1
hath	1
I	1
I	1
i'	1
it	2
julius	1
killed	1
killed	1
let	2
me	1
noble	2
so	2
the	1
the	2
told	2
you	2
was	1
was	2
with	2



term	doc. freq.	→	postings lists
ambitious	1	→	2
be	1	→	2
brutus	2	→	1 → 2
capitol	1	→	1
caesar	2	→	1 → 2
did	1	→	1
enact	1	→	1
hath	1	→	2
i	1	→	1
i'	1	→	1
it	1	→	2
julius	1	→	1
killed	1	→	1
let	1	→	2
me	1	→	1
noble	1	→	2
so	1	→	2
the	2	→	1 → 2
told	1	→	2
you	1	→	2
was	2	→	1 → 2
with	1	→	2

Implementation

index.c

- 문서집합을 역색인(inverted index) 구조로 변환한다.
- 입력 : 문서 집합 (r52.txt)
- 출력 : 사전 엔트리 목록 (dic.txt)
 역색인 파일 2종 (header.idx, posting.idx)

r52.txt

- Source: Reuters-21578 data set
- 한 라인에 하나의 문서가 저장되어 있음
 - 9100 라인(문서)
- Modified token 형태로 모두 변환됨
 - tokenized, 소문자로 정규화, stopword 제거, 기호 제거)
- DocId는 첫 줄부터 1, 2, 3, ...으로 간주

dic.txt

- 사전(dictionary)에 저장될 엔트리 목록
 - 25,611개
 - 정렬되어야 함
 - qsort 사용
 - 구조체 배열 정렬
 - Compare 함수 : 정렬 기준 1순위: 토큰(문자열), 2순위: DocId(정수)
- 예)
 - aaa
 - aabex
 - aac
 - aachener
 - aagiy
 - aaica
 - aaix

header.idx, posting.idx

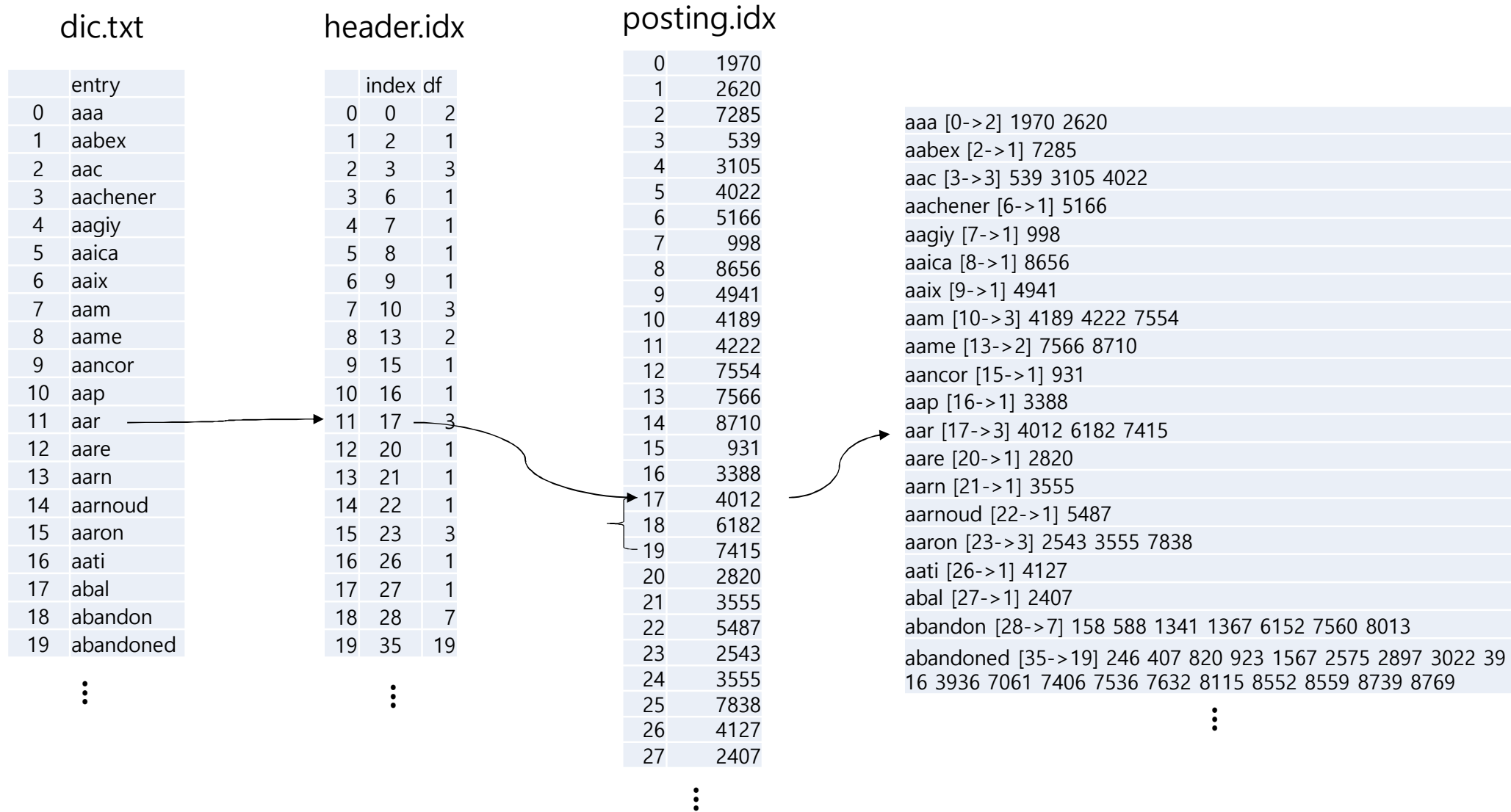
- header.idx

- 사전의 각 엔트리들의 시작 위치(index)와 문서빈도(doc. Freq.)를 저장
- 구조체 배열

```
typedef struct {  
    int index; // starting position in posting.idx  
    int df; // document frequency  
} tHEADER;
```

- posting.idx

- 문서 포스팅 정보 저장
- 정수(int) 배열



search.c

- 입력된 사용자 query를 포함하는 문서를 출력하는 프로그램
 - Query
 - exact match : 문자열
 - wildcard search : *을 포함하는 문자열
 - 불린(Boolean) 연산자 : & (and), | (or)
 - 사전 엔트리 목록 → trie 저장
 - 사용자 query → trie에서 탐색 → (역색인 구조에서) 문서 검색 → 결과 출력

프로그래밍을 위한 tip

- 구조체 정렬
 - `qsort(tokens, num_tokens, sizeof(tTokenDoc), _compare);`
- 파일 쓰기
 - `fwrite(header, sizeof(tHEADER), num, fp);`
 - `fwrite(posting, sizeof(int), num, fp);`
- 파일 읽기
 - `fread(header, sizeof(tHEADER), num, fp);`
- 동적메모리 할당/재할당
 - `malloc, realloc`
- 파일 위치 설정 및 확인
 - `fseek, ftell`

프로그래밍을 위한 tip

- trie 구조체 변경
 - 엔트리 문자열 → 엔트리 인덱스
 - 초기값은 -1로 설정

```
typedef struct trieNode {  
    int    index; // 0, 1, 2, ...  
    struct trieNode *subtrees[MAX_DEGREE];  
} TRIE;
```

- Makefile 사용
 - 복수의 소스코드 컴파일