

## Laborator nr. 3

# Construirea indecșilor direcți/indirecți

an univ. 2018 – 2019

## 1 Aplicație propusă

**1. Construiți index-ul direct cantitativ** pentru un set de documente HTML alese de voi. Indexul va stoca atât cheile (documentele indexate), cât și valorile asociate (cuvintele din cadrul documentelor) cheilor.

Aplicația va primi ca *intrare* calea către un director ce conține documentele de indexat (**Atenție:** directorul poate conține și sub-directoare). *Ieșirea* aplicației constă din **2** fișiere, astfel:

1. un prim fișier text de *mapare*: pentru fiecare document HTML (numele fișierului va fi stocat prin cale completă) se va indica numele fișierului care conține indexul direct asociat;
2. unul sau mai multe fișiere ce vor conține indexul direct determinat.

**2.** Pornind de la pseudocodul de mai jos, implementați o aplicație care va determina indexul invers corespunzător pentru colecția de documente pentru care a fost determinat indexul direct cantitativ în cadrul modulului anterior.

---

**Algorithm 1** BSBI() (preluare din [1], cap. 4)

---

```
1:  $i \leftarrow 0$ 
2: while all documents have not been processed do
3:    $n \leftarrow n + 1$ 
4:   block  $\leftarrow$  PARSENEXTBLOCK()
5:   BSBI-INVERT(block)
6:   WRITEBLOCKTODISK(block, fn)
7: end while
8: MERGEBLOCKS( f1, . . . , fn; fmerged)
```

---

Aplicația va primi ca intrare **indexul direct cantitativ** determinat în etapa anterioară, și va furniza, ca ieșire, **2** fișiere, astfel:

1. un prim fișier text de *mapare*: pentru fiecare cuvânt determinat se va indica numele fișierului care conține indexul invers asociat;
2. un fișier ce conține **indexul invers cantitativ**.

## Bibliografie

- [1] Christopher D. Manning et. al. Introduction to Information Retrieval. <http://nlp.stanford.edu/IR-book/information-retrieval-book.html>, 2009.