

# Laborator nr. 1

## Parser HTML. Extragere de cuvinte

an univ. 2018 – 2019

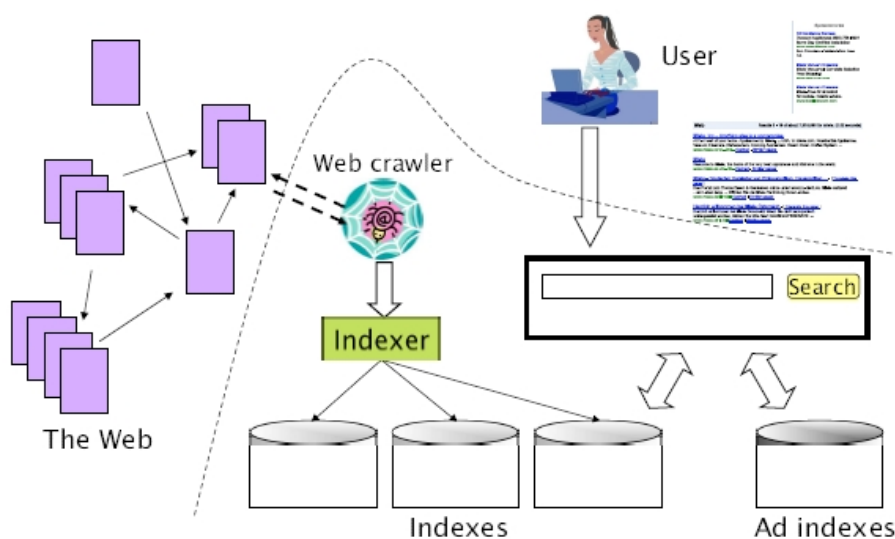


Figura 1: Structura generală a motoarelor de căutare pe WEB (preluare din [1])

## 1 Noțiuni teoretice

Caracteristici HTML:

- HTML (**H**yper**T**ext **M**arkup **L**anguage) — reprezintă un mijloc prin intermediul căruia se pot crea documente structurate;
- Autor: **Tim Berners-Lee**
- Tip limbaj: **limbaj de tip markup** — marcarea unui document pentru a stabili structura logică a acestuia
- Tip media: **text/html**
- Tip de aplicație țintă: **interpretor**
- Extensie fișiere: htm, html
- Specificații: set de elemente (etichete, tag-uri), precum și regulile de scriere ale acestora, destinate formătărilor documentelor.
- Sintaxa generală:  
`<nume-element [lista atribut]>continut</nume-element>`  
unde lista de atribute se definește ca:  
`nume-atribut="valoare-atribut" [nume-atribut="valoare-atribut"] ...`

## 2 Aplicație propusă

1. Implementați un modul de aplicație care să parseze un document html și să extragă următoarele date:

1. conținutul elementului <TITLE>, dacă acesta există;
2. conținutul atributului content, tag-ul <META>, pentru cazul în care atributul name al aceluiași element are valoarea "keywords";
3. conținutul atributului content, tag-ul <META>, pentru cazul în care atributul name al aceluiași element are valoarea "description";
4. conținutul atributului content, tag-ul <META>, pentru cazul în care atributul name al aceluiași element are valoarea "robots";
5. conținutul elementului <A>, atributul href, dacă acesta există în cadrul documentului și dacă referințele indicate NU reprezintă legături interne; link-urile trebuie reconstruite sub formă de URL-uri absolute dacă este cazul;
6. conținutul text al documentului HTML.

**Indicație:** pentru simplitate, se recomandă utilizarea librăriei JSOUP [2].

2. Realizați un modul de aplicație care primește ca intrare un text și împarte acest text în cuvinte. Pentru fiecare cuvânt în parte, se va contoriza și numărul de apariții ale celui cuvânt în textul de intrare. Rezultatele vor fi stocate fie într-un HashTable [4], fie într-un HashMap [3] (key va fi un cuvânt extras din text, value va stoca, sub forma unui întreg, numărul de apariții).

Se va considera cuvânt orice succesiune de caractere alfa-numerice cuprinsă între separatorii clasici precum: ' ' (spațiu, succesiune de spații, tab), '"', ', ', '. ', '! ', '?', etc.. În mod ideal, implementările oferite NU ar trebui să conțină bucle îmbricate.

### Observație 1:

Implementările care se bazează pe metoda split [5] a clasei String vor fi punctate cu nota maximă **8 (opt)**.

## Bibliografie

- [1] Christopher D. Manning, Prabhakar Raghavan and Hinrich Schutze. *Introduction to Information Retrieval*. Cambridge University Press, 2009.
- [2] Jonathan Hedley. jsoup: Java HTML Parser. <http://jsoup.org>.
- [3] Oracle. Class HashMap<K,V>. <http://docs.oracle.com/javase/7/docs/api/java/util/HashMap.html>.
- [4] Oracle. Class Hashtable<K,V>. <http://docs.oracle.com/javase/7/docs/api/java/util/Hashtable.html>.
- [5] Oracle. Class String. <http://docs.oracle.com/javase/7/docs/api/java/lang/String.html>.