

Proiect 2

Crawler WEB

an univ. 2018 – 2019

Aplicația corespunzătoare celei de a doua componente de proiect trebuie să implementeze un **Crawler WEB**. Acest modul trebuie să realizeze corect cereri HTTP utilizând versiunea 1.1 a protocolului și să salveze conținutul HTML al resursei indicate. Aplicația trebuie proiectată și implementată astfel încât să permită rularea în continuu. În acest sens, modulul *URL Frontier* (coada de explorare) aferent va include inițial un set de două/trei URL-uri. Rularea în continuu implică următorii pași:

1. se va prelua următorul URL din coada de explorare și se procesează astfel încât să se extragă numele domeniului explorat și URL-ul relativ al resursei dorite;
2. dacă domeniul este la **prima explorare**, atunci se va solicita resursa `/robots.txt`; dacă aceasta există, se trece la pasul 3, dacă nu se continuă cu pasul 4;
3. (dacă există `robots.txt`) se verifică clauzele `Disallow` pentru URL-ul relativ curent; dacă REP permite accesul pe resursă, atunci se trece la pasul 4, dacă nu, se trece la următorul URL din coada de explorare;
4. se preia resursa indicată de URL și se salvează local – pentru fiecare domeniu se va crea un director, apoi, în cadrul acestui director, se va urma structura de directoare din cadrul URL-ului;
5. (dacă este cazul) dacă se primește un cod 301 Moved Permanently, atunci se va reface cererea pentru noua locație și se vor actualiza datele deja salvate; orice alt tip de redirect va implica numai reluarea cererii pe noua locație a resursei, fără alte actualizări de date;
6. se analizează tag-ul HTML meta, `name="robots"`; în cazul în care este permisă extragerea link-urilor incluse în document, se vor extrage aceste link-uri sub forma unui set de URL-uri absolute; din cadrul acestui set se elimină link-urile care nu respectă REP sau care se află deja în coada de explorare;
7. se reia pasul 1.

Cerință bonus

Aplicația poate respecta următoarele cerințe bonus:

1. **Cache DNS**: aplicația implementează propriul mecanism de *caching* al înregistrărilor DNS.
2. **performanțe**: în mod secvențial, în cazul în care REP nu impune restricții de viteză, se dorește o rată medie de transfer de aproximativ 100 pagini/minut;
3. Crawler-ul WEB va fi **distribuit/paralelizat** pentru a spori performanțele; dacă, de exemplu, se vor utiliza două module de tip *fetcher*, atunci rata medie de transfer dorită va fi de cel puțin 200 pagini/minut.

Predarea celei de a doua componente de proiect se va realiza prin intermediul platformei MOODLE, similar primei componente.

Barem evaluare proiect

Proiectele pornesc ca bază de notare de la **1 punct**.

Criteriu	Punctaj
1. Realizarea corectă a cererii pentru a prelua o resursă HTML - <i>obligatoriu</i> componentă DNS – client „3rd party” <i>sau</i> componentă DNS – implementare proprie componentă HTTP – client „3rd party” <i>sau</i> componentă HTTP – implementare proprie	1 punct 2 puncte 1 punct 2 puncte
2. Salvarea completă și corectă a paginii HTML în cadrul unui singur director de lucru, fără a ține cont de structura URL-urilor în cadrul unei structuri de directoare, ținând cont de structura URL-urilor	0,5 puncte 1 punct
3. Respectarea pseudo-protocolului REP - <i>obligatoriu</i> la nivel de pagină la nivel de domeniu	1 punct 1 punct
4. Gestionarea corectă a structurilor de tip URL Frontier - <i>obligatoriu</i> reconstruirea link-urile care trebuie vizitate marcarea corectă a link-urilor care pot fi explorate marcarea corectă a link-urilor care au fost deja explorate	2 puncte

Observații

Studentii vor fi rugați să completeze o auto-evaluare a proiectului dezvoltat. Auto-evaluările care nu diferă de nota acordată de titularii de laborator vor primi **suplimentar 0.5 puncte**. În plus, implementările pot atrage punctaj suplimentar (puncte bonus) pentru media finală a disciplinei, astfel:

Implementare și gestiunea cache DNS – 2 puncte;

Rată de transfer secvențială – 100 pag./minut – 1 punct;

Paralelizare/distribuire eficientă – 2 puncte.