

# Používateľská príručka

[WebExtraction]

System pre extrakciu dát z webu

<b>Tím:</b>	č. 16, WebX
<b>Vedúci tímu:</b>	Ivan Srba
<b>Členovia tímu:</b>	Ján Brechtl, Tomáš Juhaniak, Martin Kalužník, Rastislav Krchňavý, Michal Kren, Martin Lacek, Andrej Vaculčíak
<b>Akademický rok:</b>	2016/2017
<b>Autor:</b>	Andrej Vaculčíak
<b>Verzia číslo:</b>	1.1
<b>Dátum poslednej zmeny:</b>	15.05.2017

# Obsah

Vitajte!.....	2
Úvod .....	3
Základné rozhranie .....	4
Webové sídlo .....	4
Extension pre Google Chrome .....	5
Projekt .....	6
Vytvorenie projektu.....	6
Úprava projektu a definícia dátových polí .....	7
Skript .....	8
Vytvorenie skriptu.....	8
Úprava a možnosti skriptu.....	8
Úrovne log záznamov.....	10
Nastavenie pravidelnosti spúšťania.....	10
Extrakcie.....	11
Zobrazenie extrakcií pre skript.....	11
Akcie nad dátami.....	12
Zobrazenie log záznamov.....	12
Export dát .....	12
API .....	12
Formát CSV/XLSX .....	12
Extension.....	13
Pridanie do prehliadača .....	13
Výber elementov .....	13
Postprocessing extrahovaných dát .....	17

# 1 Vitajte!

V prvom rade sa Vám chceme v mene tímu WebX poďakovať za vybratie si nášho produktu, ktorého hlavnou úlohou je pomôcť Vám získať, pomerne rýchlo, jednoducho, no hlavne pravidelne a automaticky, pre Vás užitočné dáta vo veľkom množstve a z rôznych webových sídiel.

Želáme Vám čo najväčšie množstvo úspešne extrahovaných dát!

Tím WebX

## **2 Úvod**

Táto používateľská príručka má za úlohu objasniť základné pojmy spojené s extrakciou dát, ktoré sú použité v našom systéme. Okrem toho obsahuje opis základných rozhraní a opis potrebných krokov pre úspešnú extrakciu, spolu s ukázkami obrazoviek všetkých potrebných náležitostí.

### 3 Základné rozhranie

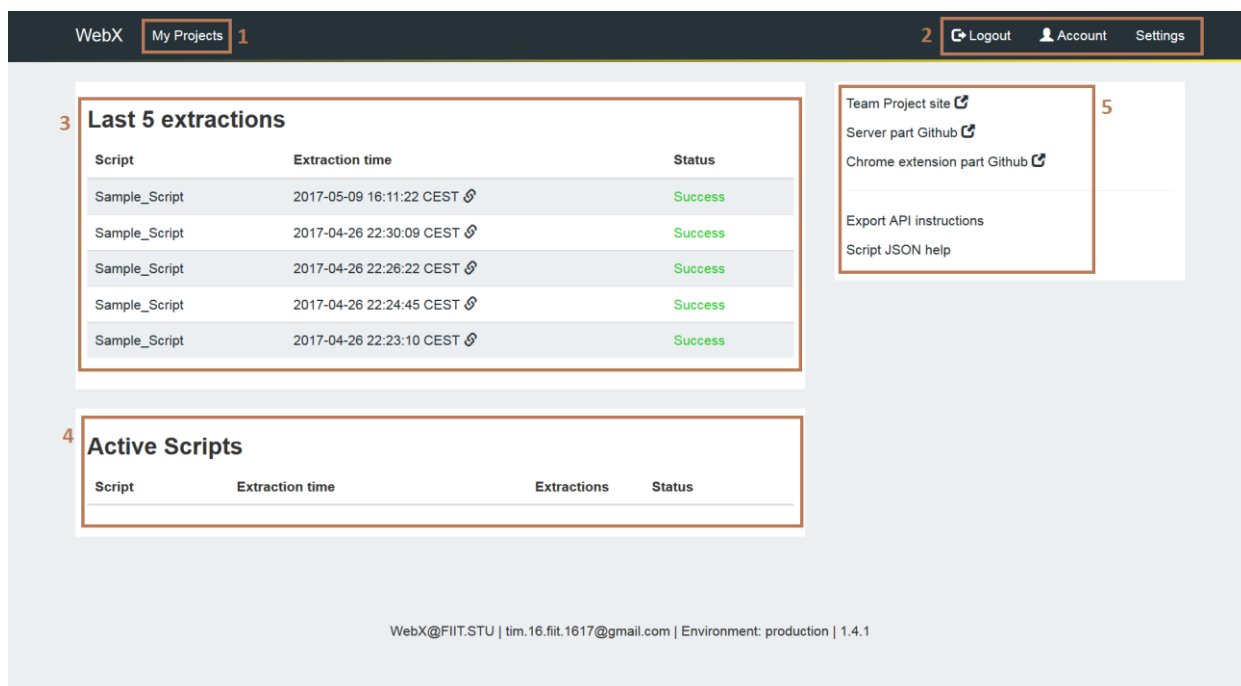
Služba sa skladá z dvoch základných častí:

1. Webové sídlo
2. Extension pre webový prehliadač

Pomocou webového sídla si používateľ spravuje svoje osobné údaje a definuje a upravuje projekty, skripty a vidí výsledky extrakcií, ktoré prebehajú alebo prebehli.

Extension pre prehliadač slúži najmä k interaktívnemu priradeniu konkrétnych polí na konkrétnych web stránkach k definovaným dátovým poliam v rámci projektu a skriptu.

#### 3.1 Webové sídlo

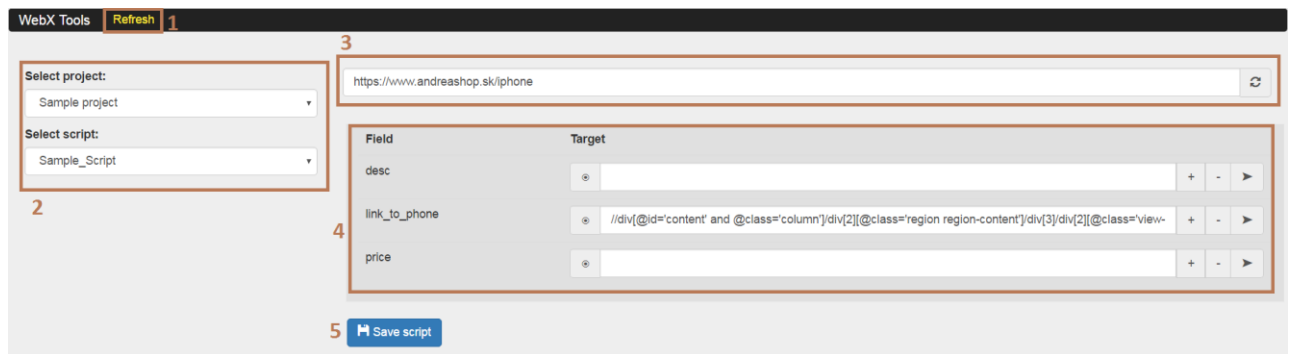


**Obr. 1 – Domovská obrazovka webového sídla po prihlásení sa**

Rozhranie domovskej stránky, zobrazenej na Obr. 1, pozostáva z týchto častí:

1. Tlačidlo na prechod do používateľom vytvorených projektov
2. Tlačidlá pre odhlásenie zobrazenie účtu a nastavení používateľových údajov
3. Tabuľka s prehľadom posledných ukončených 5 extrakcií
4. Tabuľka s prehľadom aktuálne bežiacich extrakcií
5. Externé odkazy pre stránku tímu, repozitár webového sídla + rozšírenia do prehliadača na Github-e a odkazy na stránky s inštrukciami a popisom API pre export dát a tvar JSONu, ktorý tvorí skript (viď kap. 5)

## 3.2 Extension pre Google Chrome



Obr. 2 – Základné rozhranie rozšírenia do prehliadača Google Chrome

Rozšírenie pre prehliadač Chrome má tieto základné časti:

1. Tlačidlo pre aktualizáciu dát v rozšírení (v prípade, že sme uskutočnili úpravu projektu/skriptu vo webovej časti)
2. Dropdown polia pre výber konkrétneho projektu a skriptu
3. Pole pre URL stránky z ktorej chceme dáta extrahovať
4. Zoznam dátových polí a k nim prislúchajúce XPath-y + tlačidlá pre operácie s nimi
5. Tlačidlo na uloženie aktuálneho stavu skriptu

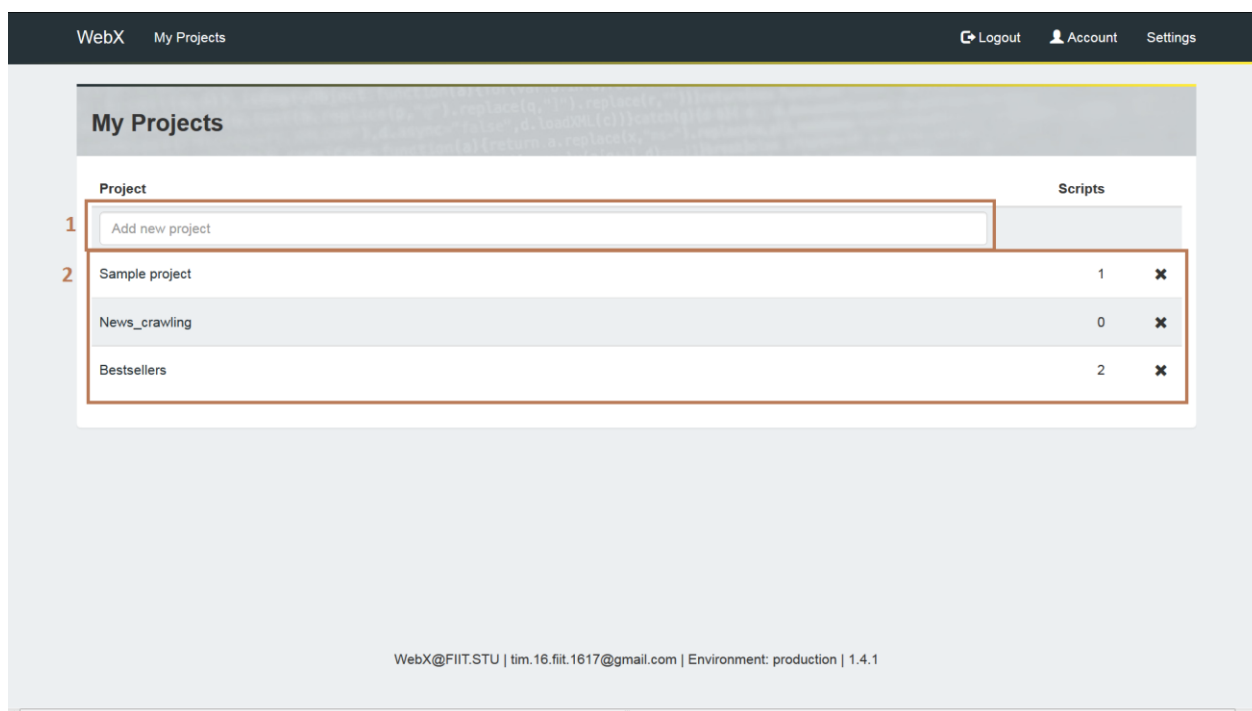
## 4 Projekt

Projekt sa v kontexte nášho systému chápe ako doména, v ktorej chceme získavať dáta. Môžeme chcieť napríklad extrahovať informácie z rôznych e-shopov o notebookoch alebo z viacerých spravodajských portálov získavať rôzne titulky článkov. Ako ďalší príklad môže byť viacero stránok s počasím, ktoré chceme zbierať a vytvoriť si napríklad podklady (datasets), na ktorých postavíme náš systém a pod.

### 4.1 Vytvorenie projektu

Ak sa chceme pracovať s projektami, stačí, keď na úvodnej stránke klikneme na tlačidlo “My Projects” (viď. Obr. 1).

Zobrazí sa nasledujúca obrazovka:

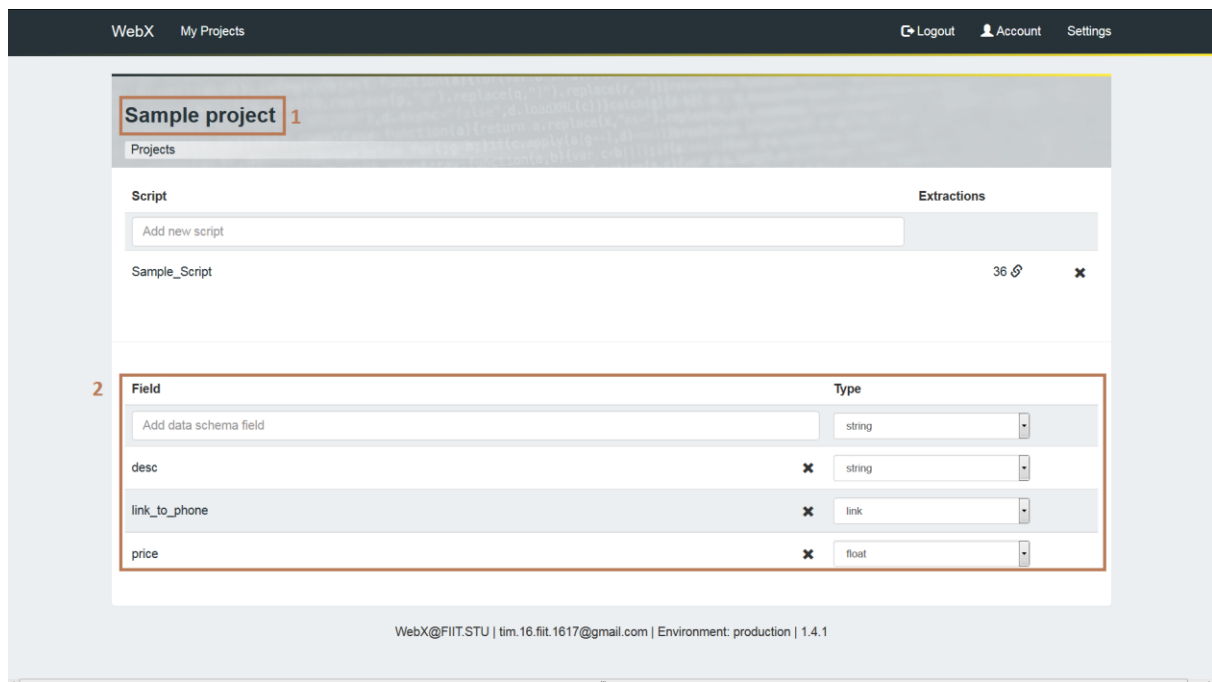


**Obr. 3 – Obrazovka pre prácu s projektami**

Na tejto obrazovke môžeme vytvoriť nový projekt (1) alebo vidíme prehľad už vytvorených projektov (2).

Do konkrétneho projektu sa dostaneme kliknutím na jeho názov.

## 4.2 Úprava projektu a definícia dátových polí



**Obr. 4 – Obrázok pre úpravu projektu a manažment dátových polí pre projekt**

Po prechode na obrázok úprav (Obr. 4), môžeme:

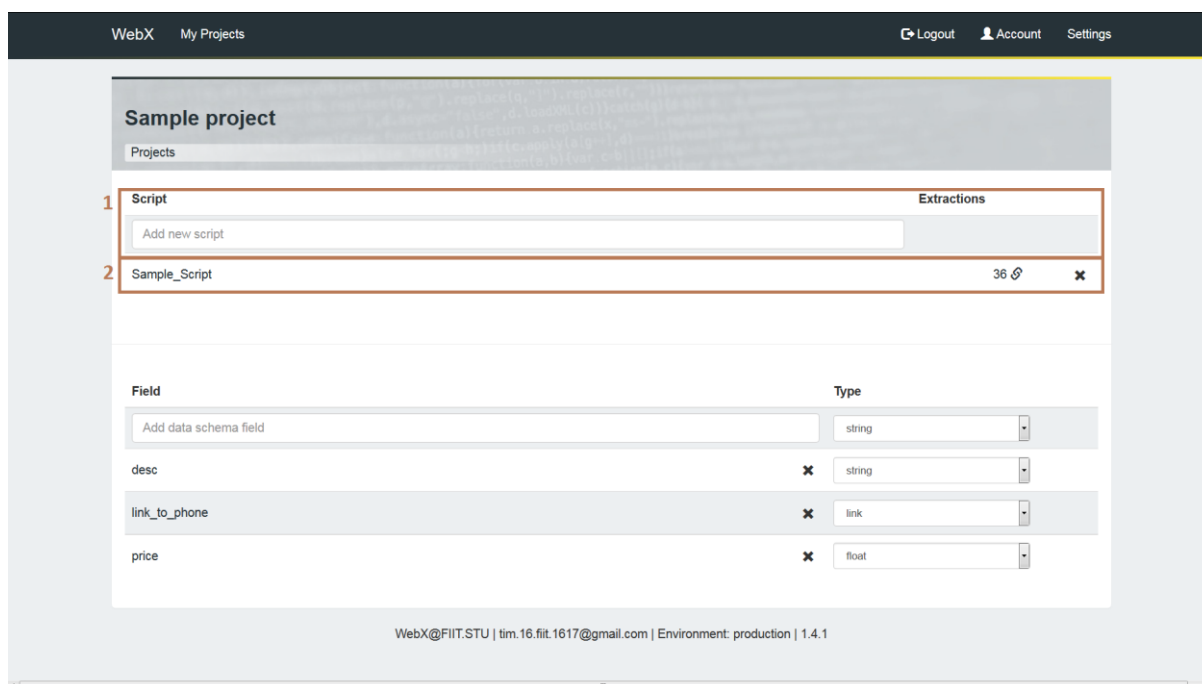
1. Upraviť názov projektu (kliknutím na názov)
2. Pridať/ upraviť existujúce dátové pole spolu s jeho typom



## 5 Skript

Každý skript predstavuje jednu web stránku v rámci danej domény, z ktorej chceme extrahovať dáta. Ak napríklad extrahujem titulky novinových článkov, vytvorím si jeden skript pre BBC News a jeden pre Al Jazeera.

### 5.1 Vytvorenie skriptu



Obr. 5 – Vytvorenie skriptu v rámci projektu

Vytvorenie nového skriptu prebieha na obrazovke správy projektu (pole 1 na Obr. 5).

Okrem toho tu vidíme aj zoznam už vytvorených skriptov (2), spolu s počtom vykonaných extrakcií.

### 5.2 Úprava a možnosti skriptu

Do úprav novo vytvoreného skriptu sa dostaneme kliknutím na jeho názov (pole 2 na Obr. 5).

Na obrazovke skriptu (Obr. 6) máme k dispozícii tieto možnosti:

1. Upraviť názov skriptu (po kliknutí na názov)
2. Upraviť JSON skriptu (v podstate teda samotný skript)

Pri úpravách skriptu je potrebné dodržiavanie určitých zásad a je potrebné poznať jeho štruktúru. Základný pomocník pre skript nájdeme na adrese:

[http://team16-16.studenti.fiit.stuba.sk/webx/json\\_help](http://team16-16.studenti.fiit.stuba.sk/webx/json_help)

3. Prehľad informácií z poslednej vykonanej extrakcie (pokiaľ nejaká prebehla) spolu s jednoduchou štatistikou prázdnych polí.

Po kliknutí na text s dátumom poslednej extrakcie sa zobrazia extrahované dáta v tejto extrakcii.

4. Nastavenie úrovne log záznamov (viac v kap. 5.2.1)
5. Uložiť úpravy JSONu a zmenu úrovne log záznamov
6. Zobrazit' zoznam extrakcií pomocou API ("API List Extractions")

Zobrazit' dáta posledných 10 extrakcií ("API Export Extractions")

Spustiť extrakciu ("Run Now")

Zobrazit' log záznamy poslednej extrakcie ("Last extraction logs")

Zobrazit' prehľad extrakcií pre daný skript ("Extractions")

7. Definovať nový interval vykonávania extrakcie (viac v kap. 5.2.2)

8. Upraviť už existujúci interval spúšťania extrakcie

**Sample\_Script 1**

Projects / Sample project

**2** Script

```
{
  "url": "https://www.andreashop.sk/iphone",
  "data": [
    {
      "name": "link_to_phone",
      "xpath": "//div[@id='content' and @class='column']/div[2]
[@class='region region-content']/div[3]/div[2][@class='view-content']
/article/div[1][@class='top']/h2/a",
      "postprocessing": [
        {
          "type": "nested",
          "data": [
            {
              "name": "price",
              "xpath": "//div[@id='content' and
@class='column']/div[2][@class='region region-content']/article/div[1]
[@class='node-content top']/div[2][@class='product-box
box-cart']/div[4][@class='product-box-row
row-cart']/div[@class='ceny']/div[1][@class='cena']/div[2]

```

**3** Last extraction: 2017-05-09 16:11:22 CEST

Execution time	1m 42s
Instances	30
Empty fields	5
Status	Success
Data Field	Empty fields
desc	5
link_to_phone	0
price	0

**4** Log level: debug

**5** Save

**6** API List Extractions API Export Extractions Run Now Last extraction logs Extractions

**7** Interval management table:

Interval	Period	First execution
5	hour	2017-05-09 22:00

**8**

Obr. 6 – Obrazovka pre úpravy skriptu

## 5.2.1 Úrovně log záznamov

Ku každej extrakcií je možné vytvárať log záznamy rôznej úrovne. Predvolená hodnota je „Error“.

Označenie úrovne	Popis
<b>Debug</b>	Najpodrobnejšia úroveň. Uchováva podrobné informácie o stave extrakcie.
<b>Warning</b>	Menej podrobná. Uchováva najmä neočakávané udalosti a chyby počas extrakcie
<b>Error</b>	Najmenej podrobná. Uchováva len chybové hlásenia extrakcie.

## 5.2.2 Nastavenie pravidelnosti spúšťania

Pre každý skript je možné nastaviť jeden alebo viac intervalov spúšťania. Pri definovaní tohto intervalu (pole 7 na Obr. 6) sa nastavujú nasledovné parametre:

- Interval – koľkokrát uplynie perióda (Period), kým sa daná extrakcia znova spustí
- Period – časové kvantum pre jeden Interval (minúta, hodina, deň, mesiac)
- First Execution – dátum a presný čas prvej extrakcie

Príklad:

Chceme daný skript vykonávať každých 5 hodín. Hodnoty budú nasledovné:

Interval = 5, Period = hour

Ak chceme vykonávať skript každých 15 minút, bude situácia nasledovná:

Interval = 15, Period = minute

## 6 Extrakcie

Extrakcia je vo svojej podstate samotná operácia, počas ktorej sa podľa JSONu v skripte extrahujú konkrétne dáta určené xPathom pomocou rozšírenia do prehliadača Chrome. Spúšťa sa pravidelne, automaticky, podľa používateľom definovaného intervalu (kap. 5.2.2) alebo manuálne po kliknutí na tlačidlo na obrazovke pre prehľad a úpravy konkrétneho skriptu (pole 6 na Obr. 6).

Každá extrakcia má svoje špecifické hodnoty atribútov. Medzi atribúty zaraďujeme:

- Čas vykonávania – koľko daná extrakcia trvala
- Status – konečný stav extrakcie (Success, Running, Fail)

Okrem toho má aj odvodené atribúty:

- Počet prázdnych polí – ak sa extrahuje prázdna hodnota
- Počet inštancií – počet objektov v rámci extrakcie

### 6.1 Zobrazenie extrakcií pre skript

Executed	Execution time	Instances	Empty fields	Status	Actions	Export
2017-05-09 16:11:22 CEST	1m 42s	30	5	Success	Data   Logs	CSV   XLSX
2017-04-26 22:30:09 CEST	1m 59s	30	5	Success	Data   Logs	CSV   XLSX
2017-04-26 22:26:22 CEST	6s 870ms	0	0	Success	Data   Logs	CSV   XLSX
2017-04-26 22:24:45 CEST	8s 162ms	0	0	Success	Data   Logs	CSV   XLSX
2017-04-26 22:23:10 CEST	8s 870ms	1	1	Success	Data   Logs	CSV   XLSX
2017-04-26 22:22:09 CEST	9s 874ms	1	1	Success	Data   Logs	CSV   XLSX
2017-04-26 22:20:36 CEST	6s 599ms	1	1	Success	Data   Logs	CSV   XLSX
2017-04-26 22:19:36 CEST	3s 853ms	0	0	Success	Data   Logs	CSV   XLSX
2017-04-26 22:17:36 CEST	2s 701ms	1	0	Success	Data   Logs	CSV   XLSX
2017-04-26 22:16:25 CEST	4s 931ms	1	1	Success	Data   Logs	CSV   XLSX
2017-04-26 22:10:17 CEST	7s 086ms	0	0	Success	Data   Logs	CSV   XLSX
2017-04-26 22:03:46 CEST	4s 738ms	1	0	Success	Data   Logs	CSV   XLSX
2017-04-26 22:02:42 CEST	9s 438ms	1	0	Success	Data   Logs	CSV   XLSX

Obr. 7 – Prehľad extrakcií pre konkrétny skript

Do zobrazenia všetkých extrakcií pre konkrétny skript sa dostaneme tlačidlom „Extractions“ (pole 6 na Obr. 6). Nájde tam prehľad základných informácií o každej extrakcii (pole 1) spolu s dostupnými akciami (pole 2).

## **6.2 Akcie nad dátami**

Pre každú extrakciu vieme buď zobraziť konkrétne extrahované dáta alebo zobraziť log záznamy pre konkrétnu extrakciu (pole 2 na Obr. 7).

### **6.2.1 Zobrazenie log záznamov**

Pri zobrazovaní log záznamov je možné tieto záznamy filtrovať podľa úrovne (viď kap. 5.2.1)

### **6.2.2 Export dát**

Všetky dáta, ktoré získame pomocou extrakcie sú uložené na serveri, z ktorého si ich vieme stiahnuť buď pomocou API alebo v štandardizovanej forme CSV, prípadne XLSX.

### **6.2.3 API**

Export dát pomocou API je možný buď priamo pomocou tlačidiel (pole 6 na Obr. 6) alebo priamo dopytom, pričom je ale potrebné vedieť tento dopyt korektne odoslať. Pomôcku nájdeme na adrese: <http://team16-16.studenti.fiit.stuba.sk/webx/export-api-instructions>

### **6.2.4 Formát CSV/XLSX**

Pokiaľ chceme dáta stiahnuť v prehľadnej forme, môžeme použiť tlačidlá na export vo formáte CSV alebo XLSX. Tie sú dostupné na viacerých miestach.

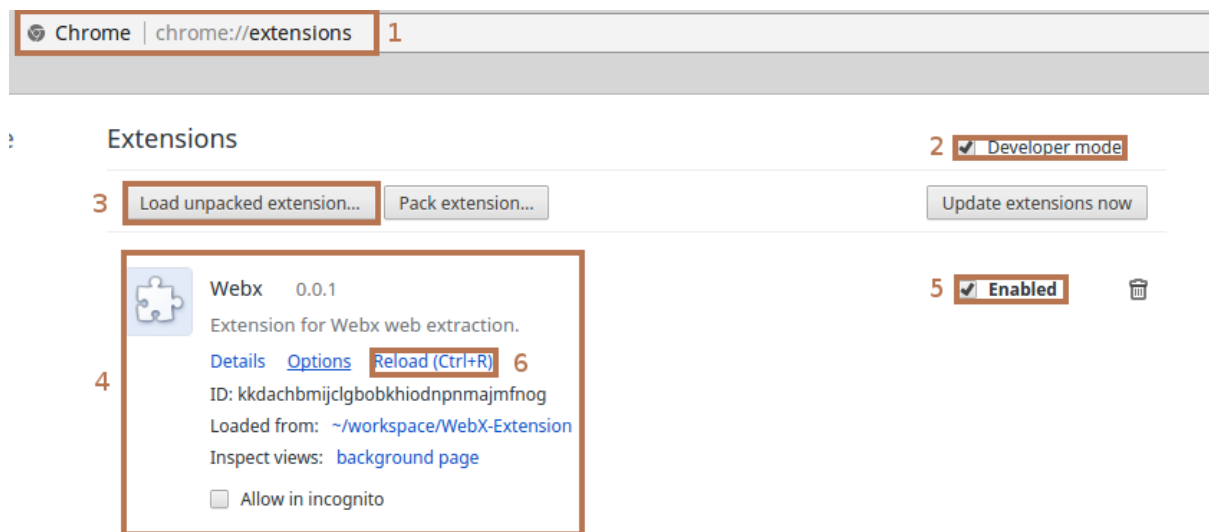
## 7 Extension

Extension je rozšírenie do Chrome prehliadača, pomocou ktorého sa vytvára skript (vo formáte JSON) pre projekt. Na základe skriptu sú následne sťahované dáta zo stránky. Najprv je potrebné si stiahnuť extension z repozitáru, ktorý sa nachádza na adrese: <https://github.com/michal55/WebX-Extension>

### 7.1 Pridanie do prehliadača

Na to, aby bolo možné používať rozšírenie v Chrome prehliadači, je potrebné najprv toto rozšírenie pridať do Chrome prehliadača. Na obrázku 8 sú naznačené nasledovné kroky:

1. V novom okne najprv načítať link `chrome://extensions/`.
2. Povolit' developer mode.
3. Načítanie samotného extensionu z priečinku, v ktorom je uložené na disku.
4. Po správnom načítaní sa rozšírenie zobrazí na obrazovke.
5. Možnosť deaktivovať alebo aktivovať rozšírenie.
6. Možnosť opakovane načítať rozšírenie v prípade, že nastali nejaké zmeny v rozšírení.



Obr. 8 – Pridanie extensionu do Chrome prehliadača

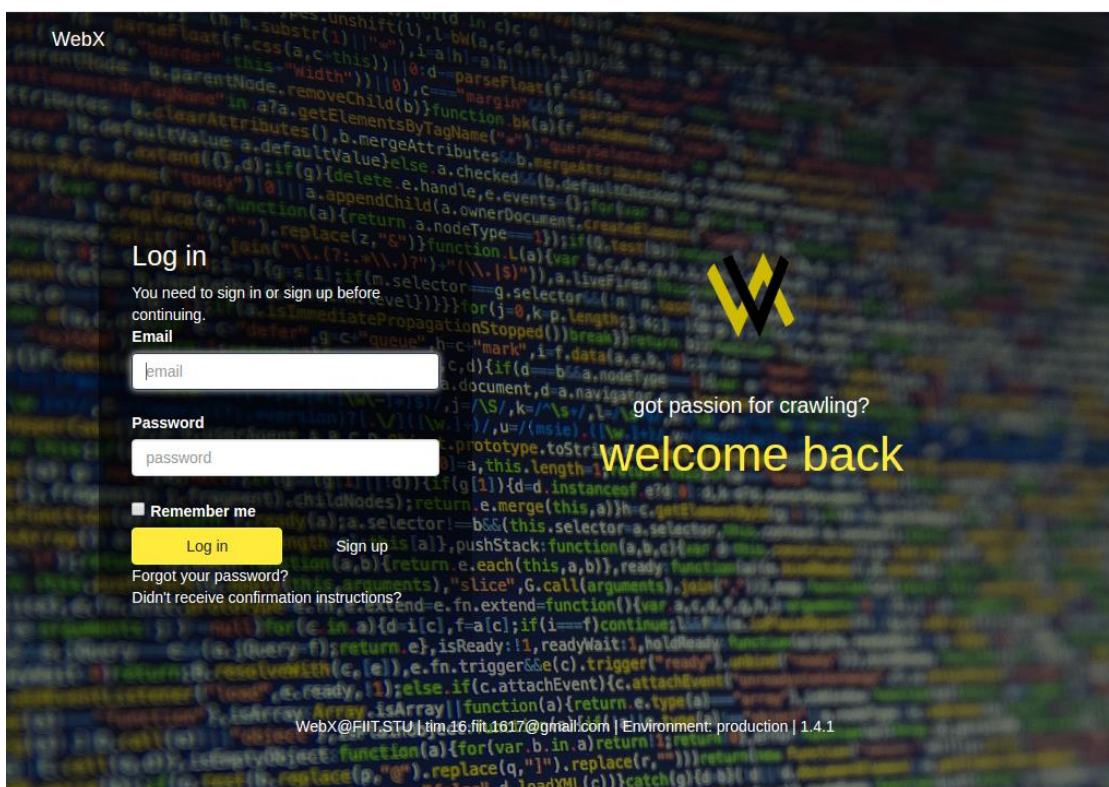
### 7.2 Výber elementov

Po načítaní extensionu do Chrome prehliadača je možné ho začať používať. Najprv je potrebné otvoriť developer tools (Vývojárske nástroje), štandardne je to za pomoci klávesy F12. Obrazovka na obrázku 9 zobrazuje extension pred prihlásením:

1. Extension medzi možnosťami vývojárskych nástrojov.
2. Tlačidlo pre prihlásenie. Po stlačení sa zobrazí obrazovka pre prihlásenie sa do webovej aplikácie - obrázok 10.
3. Tlačidlo pre znovu načítanie extensionu.



Obr. 9 – Prihlásenie do extensionu



Obr. 10 – Prihlasovacie okno

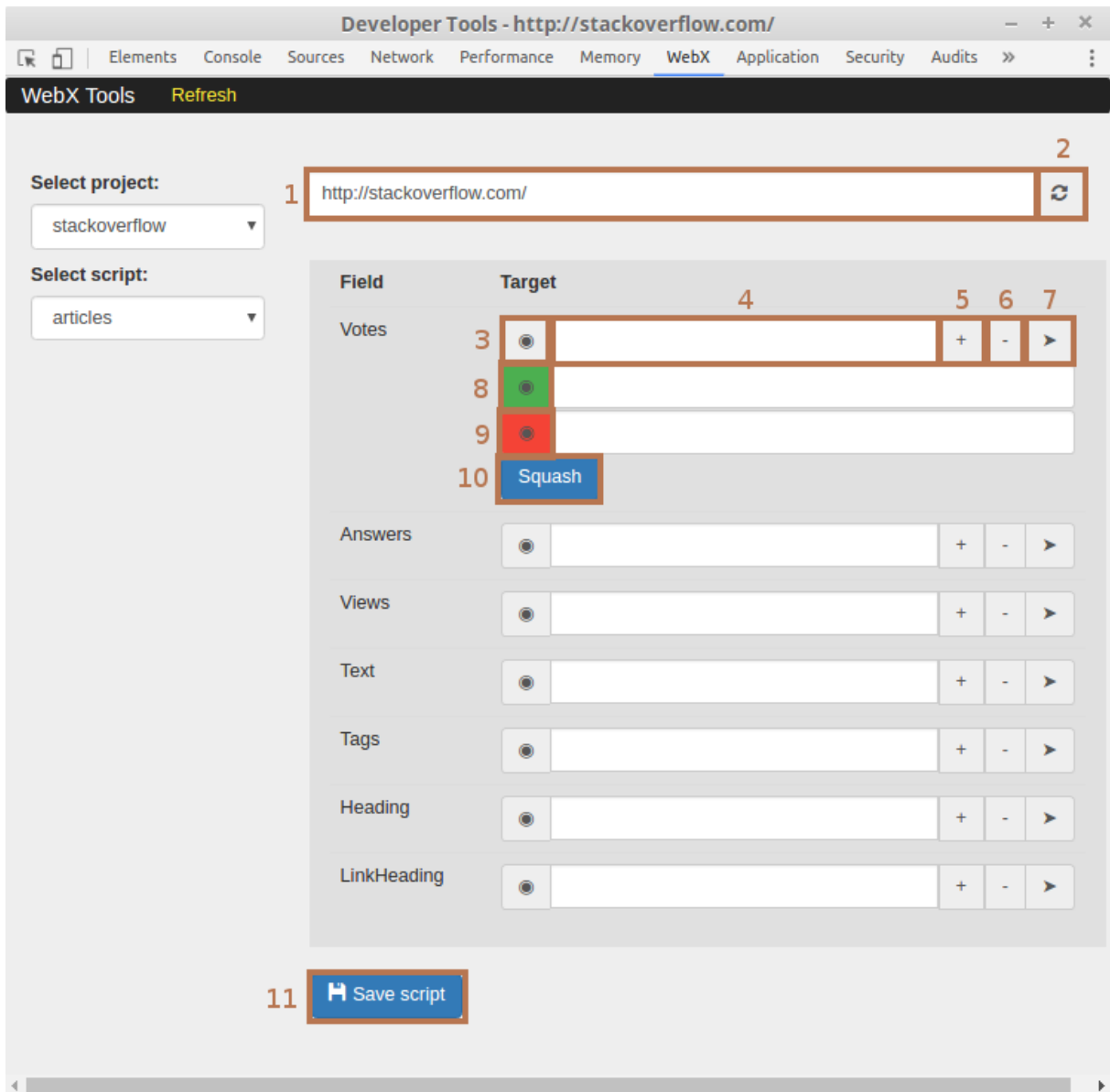
Po prihlásení sa na obrazovke vľavo zobrazí možnosť zvoliť projekt. Projekt musí už existovať vo webovej aplikácii pred tým ako je možné ho zvoliť. Po zvolení projektu sa zobrazí možnosť zvoliť konkrétny skript. Skript musí taktiež existovať vo webovej aplikácii pre zvolený projekt (v prípade vytvorenia projektu alebo skriptu po prihlásení sa do extension je možné použiť tlačidlo „Refresh“). Po zvolení projektu aj skriptu sa načítajú dátové polia (field) a k nim príslušné akcie. Na obrázku 11 je obrazovka s načítanými dátovými poľami.

1. Link stránky na ktorej skript začína (url).
2. Tlačidlo pre načítanie adresy stránky aktuálneho okna.
3. Tlačidlo na výber elementu zo stránky. Po stlačení je možné na aktuálnej stránke vybrať element. Jednotlivé elementy sa zvýrazňujú tak, ako je to zobrazené na

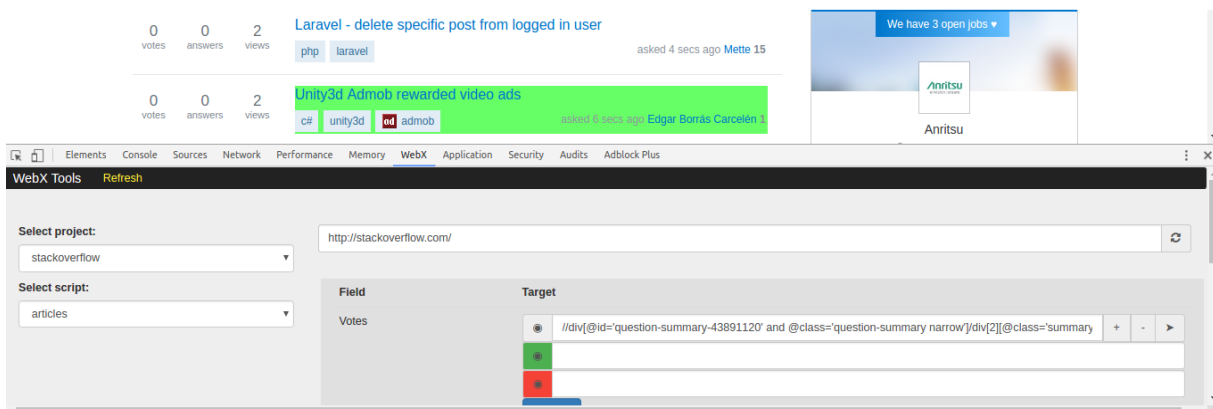
obrázku 12. Po kliknutí na konkrétny element sa načíta príslušný XPath, ktorým je daný element opísaný.

4. Pole do ktorého sa načíta XPath. V prípade potreby je možné tento XPath upravovať. Po kliknutí do poľa sa zvýrazní element, ktorý XPath opisuje.
5. Tlačidlo pre pridanie ďalšieho poľa pre XPath ku konkrétnemu dátovému poľu. Používa sa keď chceme zdefinovať XPath, ktorý opisuje pole elementov. Používa sa v kombinácii s tlačidlom z bodu 10.
6. Tlačidlo pre pridanie ďalšieho poľa pre XPath ku konkrétnemu dátovému poľu. Používa sa keď chceme zdefinovať XPath, ktorý opisuje pole elementov avšak chceme vylúčiť konkrétny element z tohto poľa. Používa sa v kombinácii s tlačidlom z bodu 10.
7. Tlačidlo pre pridanie ďalšieho spracovania elementov. Toto je opísane v časti 7.3
8. Tlačidlo s rovnakou funkcionalitou ako v bode 3.
9. Tlačidlo s rovnakou funkcionalitou ako v bode 3.
10. Tlačidlo na zjednotenie XPathov pre konkrétne dátové pole. Po stlačení sa všetky XPath-y spoja do jedného, ktorý všeobecne opisuje všetky elementy, ktoré boli zdefinované jednotlivými XPathmi.
11. Tlačidlo na uloženie skriptu.





Obr. 11 – Dátové polia v extension

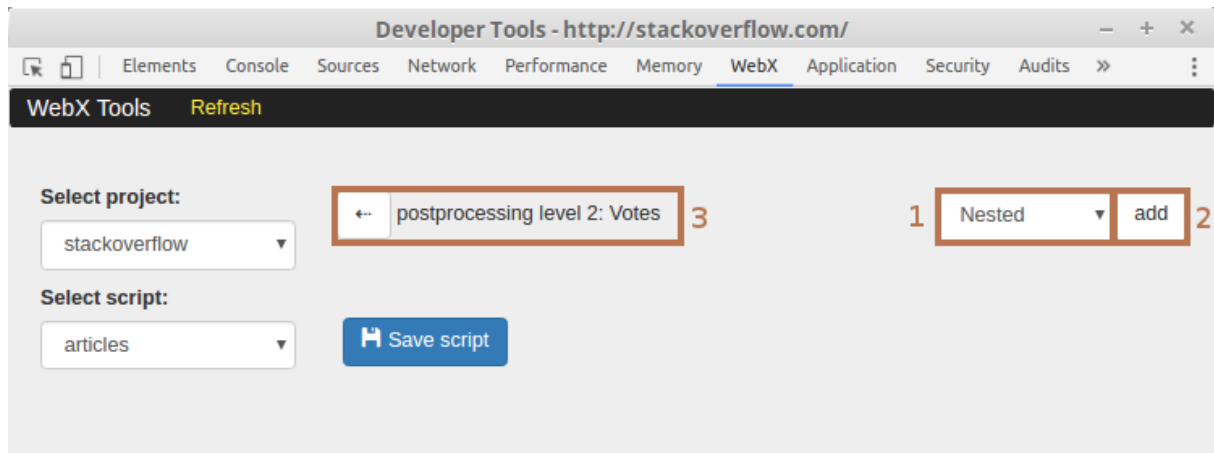


Obr. 12 – Zvýrazňovanie elementov na stránke

### 7.3 Postprocessing extrahovaných dát

Postprocessing je ďalšie spracovanie dátového poľa po extrakcii. V rozšírení je možné ho pridať cez tlačidlo opísane v časti 7.2 bodu 7. Po kliknutí na tlačidlo sa zobrazí obrazovka z obrázku 13. Na tejto obrazovke pribudli tri nové možnosti:

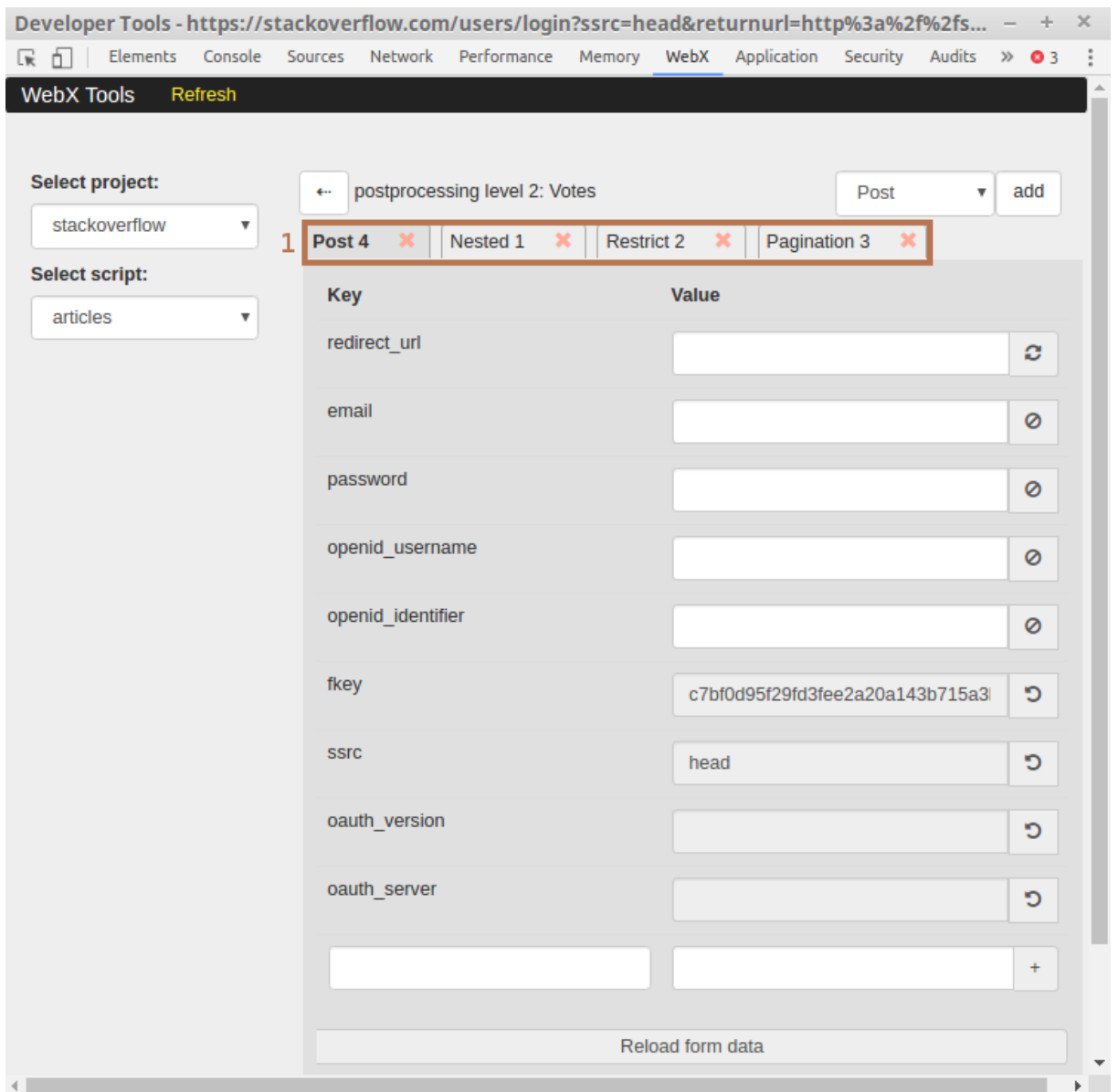
1. Výber z možností ďalšieho spracovania.
2. Pridanie konkrétneho ďalšieho spracovania.
3. Návrat na predošlú obrazovku.



Obr. 13 – Postprocessing základná obrazovka

Pre každé dátové pole je možné vybrať viacero spracovaní. Tieto sa vykonávajú v poradí, v akom boli uložené. Na obrázku 14 je obrazovka s viacerými spracovaniami. Každé spracovanie má svoje číslo a spracovania je možné medzi sebou ľubovoľne vymieňať potiahnutím myšou, alebo ich vymazať. Poradie spracovaní v skripte je postupne zľava doprava. Pridať je možné nasledovné spracovania:

Názov spracovania	Funkcionalita
Nested	Vnorenie skriptu. Používa sa keď zdefinujeme xpath, ktorý ukazuje na link, na ktorom chceme extrahovať ďalšie dáta.
Restrict	Obmedzenie na extrahovanie dát v rámci určitého elementu.
Pagination	Skript prechádza stránkovanie v prípade že je to potrebné. Používa sa na stránkach, pri ktorých sú dáta stránkované.
Attribute	Extrahovanie určitého atribútu z html prvku (elementu).
Post	Slúži na prihlásenie, ak nie je možné extrahovať dáta bez prihlásenia.
Filter	Slúži na filtrovanie dát podľa dátumu, extrahovaného v rámci dát, napríklad keď chceme dáta, ktoré boli pridané predošlý deň. (iba pre dátový typ dátum)



Obr. 14 – Příklad postprocessingov