

Export úloh 09. 11. 2021

Obsah

IN PROGRESS	4
[AMS-60] Pamäťová zložitosť elastic search na vzorke dát Created: 01/Nov/21	4
[AMS-71] Napĺňanie MongoDB priamo pri scrapovaní Created: 02/Nov/21	6
[AMS-11] Lofi používateľského rozhrania Created: 05/Oct/21 Updated: 12/Oct/21	7
[AMS-36] User databáza + zbieranie aktivity používateľa Created: 13/Oct/21 Updated: 17/Oct/21	8
TO DO	9
[AMS-95] Implementácia dummy Flask API Created: 09/Nov/21	9
[AMS-94] Návrh Flask API Created: 09/Nov/21	10
[AMS-93] Analýza možností implementácie testovania Created: 09/Nov/21	11
[AMS-92] Zabezpečenie komunikácie medzi klientom a serverom Created: 09/Nov/21	12
[AMS-91] Úprava dokumentácie zo stretnutí a šprintov tak, aby mohli ísť na stránku tímu Created: 09/Nov/21	13
[AMS-90] Integrácia MongoDB a Elasticsearch Created: 09/Nov/21	14
[AMS-89] Integrácia scrapera a vylepšeného parsera Created: 09/Nov/21	15
[AMS-88] Rozšíriť úložisko pomocou nepriradeného disku Created: 09/Nov/21	16
[AMS-87] Zabrániť zacykleným buildom Created: 07/Nov/21	17
[AMS-67] Pozrieť sa na zahodené články, kt. neobsahovali paragrafy Created: 02/Nov/21	18
[AMS-68] Analýza médií, ktoré dostaneme vo výstupoch a aké search queries treba poslať, aby sme mali výsledky ako z Google Created: 02/Nov/21	19
[AMS-72] Vytvorenie developerského prostredia Created: 02/Nov/21	20
[AMS-86] PDF report z nájdených článkov Created: 02/Nov/21	21
[AMS-85] Prihlásenie sa na odber nejakého vyhládaného človeka Created: 02/Nov/21	22
[AMS-84] Označenie článku ako videného (prečítaného) Created: 02/Nov/21	23
[AMS-83] Hodnotenie risku pri nájdených článkoch Created: 02/Nov/21	24
[AMS-82] Zvýraznenie hľadaných slovíčok Created: 02/Nov/21	25
[AMS-77] Filter na web stránke Created: 02/Nov/21	26
[AMS-46] Parsovanie tagov pomocou CSS Selector Created: 27/Oct/21	27
[AMS-42] Porovnanie výsledkov nášho vyhľadávania s výsledkami Google News Created: 13/Oct/21	28
[AMS-12] Nasu API dat na portal verejnych API v ramci propagacie Created: 05/Oct/21 ..	29
[AMS-41] Vyriešiť, čo s článkami, ktoré budú v pôvodnom zdroji nekompletné Created: 13/Oct/21	30



[AMS-37] Zistiť, ako fungujú iné zdroje Google News, prípadne aký mechanizmus používa Google News. Created: 13/Oct/21	31
[AMS-33] Zistiť, čo spája kriminálne články Created: 13/Oct/21	32
[AMS-19] Analyza - multilingvistika Created: 05/Oct/21	33
DONE	34
[AMS-57] Zavesenie klienta Created: 27/Oct/21 Resolved: 07/Nov/21	34
[AMS-59] Parsovanie tagov pomocou lxml, readability, trafiletura Created: 30/Oct/21 Resolved: 09/Nov/21	35
[AMS-73] Vyhľadanie RSS dvojčat'a k HTML článku - všeobecný postup Created: 02/Nov/21 Resolved: 09/Nov/21	38
[AMS-70] Rozšírenie google news scrapperu na lokáciu UK Created: 02/Nov/21 Resolved: 02/Nov/21	39
[AMS-51] Pridať info o projekte a opísať členov tímu Created: 27/Oct/21 Resolved: 01/Nov/21	40
[AMS-50] Request URL s odpoveďou inou ako 200 uložiť do súboru Created: 27/Oct/21 Resolved: 01/Nov/21	41
[AMS-49] Lokálna MongoDB so získanými dátami Created: 27/Oct/21 Resolved: 01/Nov/21	42
[AMS-48] Analýza bezstratovej kompresie textu Created: 27/Oct/21 Resolved: 01/Nov/21	43
[AMS-47] Parsovanie tagov pomocou regex Created: 27/Oct/21 Resolved: 30/Nov/21	44
[AMS-53] Rozšíriť docker compose Created: 27/Oct/21 Resolved: 31/Oct/21	45
[AMS-32] Prihláška na TP CUP Created: 13/Oct/21 Resolved: 02/Nov/21	46
[AMS-44] Stránka tímu Created: 18/Oct/21 Resolved: 01/Nov/21	47
[AMS-29] Získanie dát pre prototyp Created: 13/Oct/21 Resolved: 26/Oct/21	48
[AMS-43] Otvorenie portov na serveri Created: 18/Oct/21 Resolved: 18/Oct/21	49
[AMS-35] Prieskum databáz Created: 13/Oct/21 Resolved: 18/Oct/21	50
[AMS-39] Vyriešiť, aby fiitkar nemohol urobiť "sudo su" na virtuálnom stroji (aby nemal šancu sa zmeniť na roota) Created: 13/Oct/21 Resolved: 19/Oct/21	51
[AMS-31] Set up elastic search v docker container Created: 13/Oct/21 Resolved: 16/Oct/21	52
[AMS-34] Vytvoriť github projekt Created: 13/Oct/21 Resolved: 18/Oct/21	53
[AMS-13] Zistiť, ktoré miestnosti su na karticky Created: 05/Oct/21 Resolved: 18/Oct/21	54
[AMS-4] Zaobstarat stroj v škole Created: 01/Oct/21 Updated: 05/Oct/21 Resolved: 13/Oct/21	55
[AMS-10] Revizia požiadaviek Created: 05/Oct/21 Resolved: 13/Oct/21	56
[AMS-15] Realny scrapping - vytvorenie vzorky Created: 05/Oct/21 Resolved: 10/Oct/21	57

[AMS-20] Analyza parametrov Google News Created: 05/Oct/21 Resolved: 13/Oct/21 ...	59
[AMS-21] Vytvorit Slack server Created: 05/Oct/21 Resolved: 11/Oct/21	60
[AMS-3] Zoznam trestnych cinov, ktore budu pouzite ako queries na google news Created: 01/Oct/21 Resolved: 04/Oct/21	61
[AMS-5] Vysoka architektura Created: 01/Oct/21 Resolved: 03/Oct/21	62
[AMS-6] Specifikacia poziadaviek Created: 01/Oct/21 Resolved: 03/Oct/21.....	63
[AMS-7] Prieskum ziskavania dat, praca s kniznicou Created: 01/Oct/21 Resolved: 05/Oct/21	64
[AMS-8] Pripady pouzitia Created: 01/Oct/21 Resolved: 05/Oct/21	65

IN PROGRESS

[AMS-60] Pamät'ová zložitosť elastic search na vzorke dát <small>Created: 01/Nov/21</small>	
Status:	In Progress
Project:	Adverse Media Screening

Type:	Story
Assignee:	Jakub Hlavačka
Original estimate:	6 hours
Time Spent:	5 hours

Attachments:	 image-20211108-145633.png  main.py
Sprint:	AMS Sprint 1
Story point estimate:	6

Description

Zobrat dump od [Adam Šípka](#) a spravit index v elastic search, pricom jednotlivé values pre terms budu ID zo zaznamov v MongoDB

Comments

Comment by [Jakub Hlavačka](#) [07/Nov/21]

Neskôr by bolo vhodné použiť asi toto <https://hevo.com/learn/integrating-elasticsearch-and-mongodb/>

Comment by [Jakub Hlavačka](#) [08/Nov/21]

Mam zatiaľ vytvorený docker-compose, kde je mongodb a elasticsearch s ich vlastnými volumes. Mongo ma collection articles, ktorá je komprimovaná a obsahuje všetky články z articles_1.zip . K tasku príkladám script, pomocou ktorého som čítal z mongodb a snažil sa vytvárať index v elasticsearch z článkov.

Posting list sa mi zatiaľ nepodarilo naimplementovať, pretože elasticsearch analyzer ma obmedzenie pre tokenizovanie na 10000 slov...

Aby som aspoň hrubým odhadom zistil pamätovú zložitosť vytvoreného indexu v elasticsearch, tak som vkladal do dokumentu v indexe celý html body tag, pričom jeho kľúč bola id z mongodb... Prícom po nejakom 650 článku vybehol nasledovný error

Navrhujem dalsi postup:

- vyriesit problem s obmedzenim tokenizacie
- vytvarat uz rovno index (posting list), ktory budeme vyuzivat v AMS
- vytvorit dalsi task na analyzu nastroja z <https://hevo.com/learn/integrating-elasticsearch-and-mongodb/>
- Alokovať viac času !!!

EDIT: pridanie navrhu na dalsi tak s nastrojom.

[AMS-71] [Naplnanie MongoDB pri scrapovani](#) Created: 02/Nov/21

Status:	In Progress
Project:	Adverse Media Screening

Type:	Story
Assignee:	Adam Šipka
Time Spent:	1 day
Original estimate:	8 hours
Time Spent:	6 hours 30 minutes

Sprint:	AMS Sprint 1, AMS Sprint 2
Story point estimate:	8

Description

Fields: meno, link, datum vydania, region, jazyk, telo clanku

[AMS-11] [Lofi používateľského rozhrania](#) Created: 05/Oct/21 Updated: 12/Oct/21

Status:	In Progress
Project:	Adverse Media Screening

Type:	Story
Assignee:	David Silady
Time Spent:	Not Specified
Original estimate:	Not Specified

Description

Okrem Dávida bude na tomto pracovať aj Jakub Muller.

<https://www.figma.com/file/zouHpeljsUKIQPuKYDeRB2/Lo-Fi-Prototypes?node-id=0%3A1>

Poznámky: je to ok, ale story ostáva rozpracovaná, keďže sa ešte LoFi môže zmeniť príchodom rôznych funkcionalít.

[AMS-36] [User databáza + zbieranie aktivity používateľa](#) Created: 13/Oct/21 Updated: 17/Oct/21

Status:	In Progress
Project:	Adverse Media Screening

Type:	Story
Assignee:	Táňa Poláková
Time Spent:	Not Specified
Original estimate:	Not Specified

Attachments:	 Conceptual Model.bmp  Tables.bmp
---------------------	---

Comments

Comment by [Táňa Poláková](#) [17/Oct/21]

Modely su v obrazkoch.

Poznámka: nevenujeme sa tejto story zatiaľ, pretože ju momentálne nepovažujeme za až takú dôležitú.

TO DO

[AMS-95] Implementácia dummy Flask API Created: 09/Nov/21	
Status:	To Do
Project:	Adverse Media Screening
Type:	Story
Assignee:	Dominik Horvath
Original estimate:	5 hours
Issue links:	Blocks
	is blocked by AMS-94 Návrh Flask API To Do
Sprint:	AMS Sprint 2
Story point estimate:	5

[AMS-94] [Návrh Flask API](#) Created: 09/Nov/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story
Assignee:	Táňa Poláková
Original estimate:	5 hours

Issue links:	Blocks
	blocks AMS-95 Implementácia dummy Flask API To Do
Sprint:	AMS Sprint 2
Story point estimate:	5

Description

Nástroj na dokumentáciu + dokumentácia.

[AMS-93] [Analýza možností implementácie testovania](#) Created: 09/Nov/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story	Priority:	Medium
Reporter:	Táňa Poláková	Assignee:	Adam Šípka
Resolution:	Unresolved	Votes:	0
Type:	Story		
Assignee:	Adam Šípka		
Original estimate:	5 hours		

Sprint:	AMS Sprint 2
Story point estimate:	5

[AMS-92] [Zabezpečenie komunikácie medzi klientom a serverom](#) Created: 09/Nov/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story
Assignee:	David Silady
Time Spent:	Not Specified
Original estimate:	8 hours

Sprint:	AMS Sprint 05
Story point estimate:	8

Description

Obalený fetch.

[AMS-91] Úprava dokumentácie zo stretnutí a šprintov tak, aby mohli ísť na stránku tímu Created: 09/Nov/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story
Assignee:	Táňa Poláková
Time Spent:	Not Specified
Original estimate:	8 hours

Sprint:	AMS Sprint 2
Story point estimate:	8

Description

Denníky zo stretnutí, upratanie šprintov, exporthy úloh, šprint review, ...

[AMS-90] [Integrácia MongoDB a Elasticsearch](#) Created: 09/Nov/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story
Assignee:	Jakub Hlavačka
Original estimate:	8 hours

Sprint:	AMS Sprint 2
Story point estimate:	8

Description

výskum možností prepojenia, kontrola nastavení - aby neukladal elastic text, iba ID.

[AMS-89] [Integrácia scraperu a vylepšeného parseru](#) Created: 09/Nov/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story
Assignee:	Jakub Müller
Original estimate:	5 hours

Sprint:	AMS Sprint 2
Story point estimate:	5

[AMS-88] [Rozšířit úložisko pomocou nepriradeného disku](#) Created: 09/Nov/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story
Assignee:	Dominik Horvath
Original estimate:	5 hours

Sprint:	AMS Sprint 2
Story point estimate:	5

[AMS-87] [Zabrániť zacykleným buildom](#) Created: 07/Nov/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story
Assignee:	Dominik Horvath
Original estimate:	5 hours

Sprint:	AMS Sprint 2
Story point estimate:	5

Description

Github actions timeout setting aby nam nebezali buildy zbytocne dlho

[AMS-67] [Pozriet' sa na zahodené články, kt. neobsahovali paragrafy](#) Created:
02/Nov/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story
Assignee:	Unassigned

Sprint:	
----------------	--

[AMS-68] [Analýza médií, ktoré dostaneme vo výstupoch a aké search queries treba poslať, aby sme mali výsledky ako z Google](#) Created: 02/Nov/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story
Assignee:	Unassigned

Sprint:	
----------------	--

[AMS-72] [Vytvorenie developerského prostredia](#) Created: 02/Nov/21

Status:	To Do
----------------	-------

Project:	Adverse Media Screening
-----------------	---

Type:	Story
--------------	-------

Assignee:	Unassigned
------------------	------------

Sprint:	
----------------	--

[AMS-86] [PDF report z nájdených článkov](#) Created: 02/Nov/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story
Assignee:	Unassigned

Description
Zatiaľ v backlogu.

[AMS-85] [Prihlásenie sa na odber nejakého vyhľadaného človeka](#) Created: 02/Nov/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story
Assignee:	Unassigned

Sprint:	
----------------	--

Description

vždy keď sa pridá nový článok o hľadanom človeku. V backlogu.

[AMS-84] [Označenie článku ako videného \(prečítaného\)](#) Created: 02/Nov/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story
Assignee:	Unassigned

Sprint:	
----------------	--

Description	V backlogu.
--------------------	-------------

[AMS-83] [Hodnotenie risku pri najdených článkoch](#) Created: 02/Nov/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story
Assignee:	Unassigned

Sprint:	
----------------	--

Description
Používateľ si ohodnotí, aký risk má daný článok + nejaká poznámka k tomu. V backlogu.

[AMS-82] [Zvýraznenie hľadaných slovíčok](#) Created: 02/Nov/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story
Assignee:	Unassigned

Sprint:	
----------------	--

Description	V backlogu.
--------------------	-------------

[AMS-77] [Filter na web stránke](#) Created: 02/Nov/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story
Assignee:	Unassigned

Sub-tasks:	Key	Summary	Type	Status	Assignee
	AMS-78	obmedzenie vyhľadavania na základe dá...	Subtask	To Do	
	AMS-79	vyhľadanie na základe zločinov	Subtask	To Do	
	AMS-80	geografická lokácia	Subtask	To Do	
	AMS-81	boolean query	Subtask	To Do	
Sprint:					

[AMS-46] [Parsovanie tagov pomocou CSS Selector](#) Created: 27/Oct/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story
Assignee:	Jakub Müller
Original estimate:	5 hours

Sprint:	AMS Sprint 1
Story point estimate:	5

Description

Pomocou CSS Selector je potrebné vyskúšať parsovať tagy, ktoré obsahujú text

[AMS-42] Porovnanie výsledkov nášho vyhľadávania s výsledkami Google News

Created: 13/Oct/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story
Assignee:	Unassigned
Time Spent:	Not Specified
Original estimate:	Not Specified

[AMS-12] [Nasu API dat na portal verejnych API v ramci propagacie](#) Created: 05/Oct/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story
Assignee:	Unassigned
Time Spent:	Not Specified
Original estimate:	Not Specified

Description

Pridané do backlogu.

[AMS-41] Vyriešiť, čo s článkami, ktoré budú v pôvodnom zdroji nekompletné

Created: 13/Oct/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story
Assignee:	Unassigned
Time Spent:	Not Specified
Original estimate:	Not Specified

Description

Články môžu byť vymazané, prípadne iba ich obrázky. Návrhy: vytvoriť pdf, screenshot, ...

[AMS-37] Zistiť, ako fungujú iné zdroje Google News, prípadne aký mechanizmus používa Google News. Created: 13/Oct/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story
Assignee:	Adam Šípka
Time Spent:	Not Specified
Original estimate:	Not Specified

[AMS-33] [Zistiť, čo spája kriminálne články](#) Created: 13/Oct/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story
Assignee:	Unassigned
Time Spent:	Not Specified
Original estimate:	Not Specified

Description

Napríklad aké frázy a slová sa nachádzajú v dokumentoch o podozrivých osobách a v iných článkoch nie (slovo "commit").

[AMS-19] [Analyza - multilingvistika](#) Created: 05/Oct/21

Status:	To Do
Project:	Adverse Media Screening

Type:	Story		
Assignee:	Unassigned		
Σ Remaining Estimate:	Not Specified	Remaining Estimate:	Not Specified
Σ Time Spent:	Not Specified	Time Spent:	Not Specified
Σ Original Estimate:	Not Specified	Original estimate:	Not Specified

Sub-tasks:	Key	Summary	Type	Status	Assignee
	AMS-38	Zistiť, ako funguje Google Translate...	Subtask	To Do	Adam Šípka
	AMS-40	inteligentný preklad slovíčok do slov...	Subtask	To Do	
	AMS-76	Iné možnosti prekladu článkov	Subtask	To Do	

Description

API na Google translate, preklad klucovych slov. Zatiaľ v backlogu.

DONE

[AMS-57] [Zavesenie klienta](#) Created: 27/Oct/21 Resolved: 07/Nov/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	David Silady
Original estimate:	13 hours
Time spent:	13 hours, 30 minutes






Sprint:	AMS Sprint 1
Story point estimate:	13

Description
Zavesit' klienta bud' cez Express alebo priamo cez NGINX. Testovacie API calls, redux, ...

[AMS-59] [Parsovanie tagov pomocou lxml, readability, trafiletura](#) Created: 30/Oct/21
Resolved: 09/Nov/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Táňa Poláková
Original estimate:	3 hours
Time Spent:	6 hours

Attachments:	 2_sk_article.html  lxml_measurement - Copy.txt  p_lxml_measurement.txt  readability_measurement.txt  trafiletura_measurement.txt
Sprint:	AMS Sprint 1
Story point estimate:	3

Description

pozriet aj ine jazyky, pozriet dokumentaciu, preco kniznica suvisi s jazykom

Comments

Comment by [Táňa Poláková](#) [07/Nov/21]

Readability - funguje iba na zaklade nazvov tagov a tried, nie na zaklade jazyka textu - funkcia summary() extrahuje hlavny obsah HTML stranky. Ako prve vymaze vsetky script alebo style tagy. Nasledne maza tagy na zaklade regexov:

```
{ {'unlikelyCandidatesRe':  
re.compile('combx|comment|community|disqus|extra|foot|header|menu|remark|rss|shoutbox|sideb  
ar|sponsor|ad-break|agegate|pagination|pager|popup|tweet|twitter',re.I),  
'okMaybeItsACandidateRe': re.compile('and|article|body|column|main|shadow',re.I), } }  
  
if (  
    REGEXES["unlikelyCandidatesRe"].search(s)  
    and (not REGEXES["okMaybeItsACandidateRe"].search(s))  
    and elem.tag not in ["html", "body"]  
) :  
    log.debug("Removing unlikely candidate - %s" % describe(elem))  
    elem.drop_tree()
```

Potom pridava score tagom na zaklade urcilych kriterii:

```

def class_weight(self, e):
    weight = 0
    for feature in [e.get("class", None), e.get("id", None)]:
        if feature:
            if REGEXES["negativeRe"].search(feature):
                weight -= 25

            if REGEXES["positiveRe"].search(feature):
                weight += 25

            if self.positive_keywords and
self.positive_keywords.search(feature):
                weight += 25

            if self.negative_keywords and
self.negative_keywords.search(feature):
                weight -= 25

            if self.positive_keywords and self.positive_keywords.match("tag-" +
e.tag):
                weight += 25

            if self.negative_keywords and self.negative_keywords.match("tag-" +
e.tag):
                weight -= 25

    return weight
def score_node(self, elem):
    content_score = self.class_weight(elem)
    name = elem.tag.lower()
    if name in ["div", "article"]:
        content_score += 5
    elif name in ["pre", "td", "blockquote"]:
        content_score += 3
    elif name in ["address", "ol", "ul", "dl", "dd", "dt", "li", "form",
"aside"]:
        content_score -= 3
    elif name in [
        "h1",
        "h2",
        "h3",
        "h4",
        "h5",
        "h6",
        "th",
        "header",
        "footer",
        "nav",
    ]:
        content_score -= 5
    return {"content_score": content_score, "elem": elem}

```

Vyberie sa najlepší kandidát na základe score a potom sa hľadajú susedné elementy, ktoré by s najlepšími kandidátmi mohli súvisieť.

Comment by [Táňa Poláková](#) [08/Nov/21]

Readability - pri 20tich clankoch nesedel jeden extrahovany s povodnym clankom. Extrahovany hlavný obsah bol nezmysel - iba svg obrázky.

V kode nemaju velku vahu headings, co by bolo mozne upravit pre nase potreby. Vseobecne mi ale pride, ze v kode je bordel - vela zakomentovanych veci.

Comment by [Táňa Poláková](#) [08/Nov/21]

Trafilatura - vsetkych 20 clankov obsahovalo hlavný content - aj ten, ktorý readability nezvladol. Zaujímavostou bolo, ze v kode pouzili readability a na zaklade dlzky extrahovaneho textu sa rozhodli, ci pouziju svoj alebo readability algoritmus.

Rozdiel je, ze trafilatura vracia iba holy text, bez tagov.

Comment by [Táňa Poláková](#) [08/Nov/21]

Readability funguje rovnako pre anglicke aj slovenske clanky.

Priklad:

(zdroj: <https://www.cas.sk/clanok/2603841/brutalna-dvojnashobna-vrazda-v-lednickych-rovniach-kedy-sa-zacne-proces-s-obzalovanim/>)

[2_sk_article.html](#)

Trafilatura tiež.

[AMS-73] [Vyhládanie RSS dvojčat'a k HTML článku - všeobecný postup](#) Created:
02/Nov/21 Resolved: 09/Nov/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Dominik Horvath
Original estimate:	8 hours
Time Spent:	4 hours

Sprint:	AMS Sprint 1
Story point estimate:	8

[AMS-70] [Rozšírenie google news scrapperu na lokáciu UK](#) Created: 02/Nov/21 Resolved: 02/Nov/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Dominik Horvath
Original estimate:	1 hour
Time Spent:	1 hour

Sprint:	AMS Sprint 1
Story point estimate:	1

Comments

Comment by [Dominik Horvath](#) [02/Nov/21]

Scrapper je pripraveny na parsovanie clankov z lokacie UK, v jazyku anglictina. Staci inicializovat triedu na parsovanie takto:

```
gnews_parser = GnewsParser()  
gnews_parser.setup_search("covid", '2021-09-01', '2021-09-02', locale="en-gb")
```

[AMS-51] [Pridat' info o projekte a opisat' členov tímu](#) Created: 27/Oct/21 Resolved: 01/Nov/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Jakub Müller
Original estimate:	4 hours
Time Spent:	1 hour

Sprint:	AMS Sprint 1
Story point estimate:	4

[AMS-50] [Request URL s odpoved'ou inou ako 200 uložit' do súboru](#) Created:
27/Oct/21 Resolved: 01/Nov/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Adam Šípka
Original estimate:	3 hours
Time Spent:	1 hour

Sprint:	AMS Sprint 1
Story point estimate:	3



Description

Čokoľvek, čo počas scrapovania dostane odpoveď inú ako 200 je potrebné uložiť do súboru. Okrem URL treba uložiť aj zločin.

[AMS-49] [Lokálna MongoDB so získanými dátami](#) Created: 27/Oct/21 Resolved: 01/Nov/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Adam Šípka
Original estimate:	8 hours
Time Spent:	5 hours

Attachments:	 mongodb zistenia_01.txt  mongodb zistenia_02.txt
Sprint:	AMS Sprint 1
Story point estimate:	8

Description

Naše získané dáta (všetky fields) dať do lokálnej MongoDB, zistiť, aké sú možnosti práce s ňou, koľko miesta zaberá na disku - uvidíme, ako naše dáta Mongo zredukuje. Zistiť, ako je tvorené ID.

Po porovnaní nám vyšlo, že sa to oplatí viac ako metóda ZIP

[AMS-48] [Analýza bezstratovej kompresie textu](#) Created: 27/Oct/21 Resolved: 01/Nov/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Jakub Hlavačka
Original estimate:	8 hours
Time Spent:	6 hours

Sub-tasks:	Key	Summary	Type	Status	Assignee
	AMS-61	bezstratova kompresia text	Subtask	Done	
	AMS-62	analyza vyhladavania nad komprimovany...	Subtask	Done	
	AMS-63	elastic search a komprimacia	Subtask	Done	
	AMS-64	realne kniznice	Subtask	Done	
	AMS-65	vyskusat nad datach	Subtask	Done	
	AMS-66	Ako mongoDB pracuje s vopred komprimo...	Subtask	Done	
Sprint:	AMS Sprint 1				
Story point estimate:	8				

Description

Analýza bezstratovej kompresie text; analýza vyhládavania nad komprimovaným textom; elastic search a komprimácia; encoding tetxtu; reálne knižnice; vyskúšať na dátach; ako MongoDB pracuje s vopred komprimovaným textom.

Poznámka vedúceho: úloha je splnená, jej podúloha bola preradená do novej story z dôvodu, že priamo nesúvisela s touto Story. Bola testovaná metóda ZIP

[AMS-47] [Parsovanie tagov pomocou regex](#) Created: 27/Oct/21 Resolved: 30/Nov/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Táňa Poláková
Original estimate:	8 hours
Time Spent:	2 hours

Sprint:	AMS Sprint 1
Story point estimate:	8

Description

Pomocou regexov je potrebné vyskúšať rozparsovať tagy, ktoré obsahujú text. Dôležité je, aby mali správne poradie.

Comments

Comment by [Táňa Poláková](#) [30/Oct/21]

Analýza ukázala, že nie je efektívne pracovať s regexami. Úloha sa ukončila a nahradila ju nová <https://tim8-2021.atlassian.net/browse/AMS-59>

[AMS-53] [Rozšíriť docker compose](#) Created: 27/Oct/21 Resolved: 31/Oct/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Dominik Horvath
Original estimate:	13 hours
Time Spent:	7 hours

Sub-tasks:	Key	Summary	Type	Status	Assignee
	AMS-54	Kontajner Flask API	Subtask	Done	
	AMS-55	Kontajner Express server	Subtask	Done	
	AMS-56	NGINX konfigurácia	Subtask	Done	
Sprint:	AMS Sprint 1				
Story point estimate:	13				

Comments

Comment by [Dominik Horvath](#) [31/Oct/21]

Vytvorene 2 github repozitare: [node server](#) a [flask server](#). Oba repozitare obsahuju github actions procedury nakonfigurovane tak, aby sa spustili pri pushnuti / mergnuti do main branche.

Github actions vykonaju build kontajneru z aktualnej main branch, nahra hotovy image na docker-hub a nasledne nacita najnovsie docker image aj na timovy virtualny stroj, kde ich rovno spusti a nahradi stare kontajneru.

Nove cesty v nginx reverse proxy:

- /api => flask server
- /ams => node server

[AMS-32] [Prihláška na TP CUP](#) Created: 13/Oct/21 Resolved: 02/Nov/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Unassigned
Time Spent:	Not Specified
Original estimate:	Not Specified

[AMS-44] [Stránka tímu](#) Created: 18/Oct/21 Resolved: 01/Nov/21


Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Jakub Müller
Time Spent:	Not Specified
Original estimate:	Not Specified

[AMS-29] [Získanie dát pre prototyp](#) Created: 13/Oct/21 Resolved: 26/Oct/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story		
Assignee:	Dominik Horvath		
Σ Remaining Estimate:	Not Specified	Remaining Estimate:	Not Specified
Σ Time Spent:	Not Specified	Time Spent:	Not Specified
Σ Original Estimate:	Not Specified	Original estimate:	Not Specified

Attachments:	 image-20211019-151711.png				
Sub-tasks:	Key	Summary	Type	Status	Assignee
	AMS-30	Redukcia a kategorizacia zoznamu zloč...	Subtask	To Do	Dominik Horvath

Comments

Comment by Jakub Hlavačka [19/Oct/21]
https://github.com/JKBGIT1/scraper
V readme je sposob akym to spustit.
Treba osetrit veci, ktore su v komentoch spider.py ako TODO
Comment by Jakub Hlavačka [19/Oct/21]
<attachment>
Asi to ide moc rychlo...
Comment by Jakub Hlavačka [19/Oct/21]
Zjavne bude potrebne zabezpecit, aby to nerobilo requesty na server, kde uz dostal timeout...

[AMS-43] [Otvorenie portov na serveri](#) Created: 18/Oct/21 Resolved: 18/Oct/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	David Silady
Time Spent:	Not Specified
Original estimate:	Not Specified

Description

Otvorené:
443, 80, 22 (asi aj 53)
Je lepšie ostatné ani neotvárať. (Pochybujem, že by to vôbec šlo bez ďalších komplikácií - Docker)

[AMS-35] [Prieskum databáz](#) Created: 13/Oct/21 Resolved: 18/Oct/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Jakub Müller
Time Spent:	Not Specified
Original estimate:	Not Specified

Description

Všetko v containeroch. Analýza MongoDB, ... kompatibilita s elastic search

[AMS-39] Vyriešiť, aby fiitkar nemohol urobiť "sudo su" na virtuálnom stroji (aby nemal šancu sa zmeniť na roota) Created: 13/Oct/21 Resolved: 19/Oct/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Táňa Poláková
Time Spent:	Not Specified
Original estimate:	Not Specified

Description

a nech sa ani nevie prepnut ani na ubuntu

Comments

Comment by [Táňa Poláková](#) [17/Oct/21]
<https://www.thegeekdiary.com/how-to-disable-sudo-su-for-users-in-sudoers-configuration-file/>

[AMS-31] [Set up elastic search v docker container](#) Created: 13/Oct/21 Resolved: 16/Oct/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Jakub Hlavačka
Time Spent:	Not Specified
Original estimate:	Not Specified

Attachments:	 docker-compose.yml  image-20211016-144056.png
---------------------	--

Comment by [Jakub Hlavačka](#) [16/Oct/21]

<attachment>

Pre spustenie je potrebne byt v adresari /home/fiitkar/docker-folder a odpalit command: *sudo docker-compose up -d*

Prepinac -d zabezpeci, ze sa proces spusti na pozadi.

Command na zastavenie a zmazanie docker kontajneru: *sudo docker-compose down*

Nerobit: *sudo docker-compose down -v* , pretoze to zmaze volume, v ktorom su ulozene naindexovane data.

Comment by [Jakub Hlavačka](#) [16/Oct/21]

[Dominik Horvath](#) mozes indexovat data na masine.

[AMS-34] [Vytvorit' github projekt](#) Created: 13/Oct/21 Resolved: 18/Oct/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Jakub Müller
Time Spent:	Not Specified
Original estimate:	Not Specified

Description

Viac repozitárov pod jedným projektom

[AMS-13] [Zistit, ktore miestnosti su na karticky](#) Created: 05/Oct/21 Resolved: 18/Oct/21

Status:	Done
Project:	Adverse Media Screening


Type:	Story
Assignee:	David Silady
Time Spent:	Not Specified
Original estimate:	Not Specified

Description

6-te poschodie - coworking.
Treba sa 5 krát pipnuť a napísať pánovi (neviem komu).
Je tam stále restricted access - treba poznať správnych ľudí.
Vraj sa tam na tím FIIT-WIX škaredo pozerali.
Netreba tam mať rúško?

[AMS-4] [Zaobstarat stroj v škole](#) Created: 01/Oct/21 Updated: 05/Oct/21 Resolved: 13/Oct/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Táňa Poláková
Time Spent:	Not Specified
Original estimate:	Not Specified
Attachments:	 virtualne-stroje-info2021.txt

Description

- VM na fiit(zdarma)
- 1GB Ram, 1CPU, min 50 GB - 100GB disk
- nie FreeBSD, ale skor Ubuntu

Comments

Comment by [Táňa Poláková](#) [03/Oct/21]

Kontaktovala som Ing. Juraja Petríka (5731@is.stuba.sk), ktorý by podľa úvodnej prezentácie mal mať na starosti virtuálne stroje.

Comment by [Táňa Poláková](#) [05/Oct/21]

[virtualne-stroje-info2021.txt](#)

Žiadosť o SSH key bola odoslaná na meno polakova18

[AMS-10] [Revizia požiadaviek](#) Created: 05/Oct/21 Resolved: 13/Oct/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Táňa Poláková
Time Spent:	Not Specified
Original estimate:	Not Specified

Description

Zpracovať poznámky od vedúceho zo story [AMS-6].

Schválené vedúcim, ale aj tak sa môžu ešte v priebehu vývoja softvéru meniť.

[AMS-15] [Realny scrapping - vytvorenie vzorky](#) Created: 05/Oct/21 Resolved: 10/Oct/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story		
Assignee:	Dominik Horvath		
Σ Remaining Estimate:	Not Specified	Remaining Estimate:	Not Specified
Σ Time Spent:	Not Specified	Time Spent:	Not Specified
Σ Original Estimate:	Not Specified	Original estimate:	Not Specified

Sub-tasks:	Key	Summary	Type	Status	Assignee
	AMS-16	Ziskanie vzorky RSS	Subtask	Done	Jakub Hlavačka
	AMS-17	API / scrapovanie	Subtask	Done	Dominik Horvath

Description

Ziskanie vzorky: Jeffrey epstein, poslednych 5 rokov.

Subtasks:

[AMS-16] [Ziskanie vzorky RSS](#) Created: 05/Oct/21 Updated: 08/Oct/21 Resolved: 08/Oct/21

Status:	Done
Project:	Adverse Media Screening

Type:	Subtask
Assignee:	Jakub Hlavačka
Time Spent:	Not Specified
Original estimate:	Not Specified

Attachments:	google_news_rss.py
---------------------	--------------------

Description

ziskanie RSS a zistenie, ci je to efektivne alebo nie

+ casovanie vsetkeho, co sa vykona

Comments

Comment by [Jakub Hlavačka](#) [08/Oct/21]

RSS vrati 100 clankov na danu query.

Total time: 0.046875

Malo by to byt v sekundach.

Ak sa vyberieme touto cestou, tak bude zjavne potrebne sa dostat k dalsim clankom a nie len prvym 100, ktore to vyhodi...

[AMS-17] [API / scrapovanie](#) Created: 05/Oct/21 Resolved: 10/Oct/21

Status: Done

Project: [Adverse Media Screening](#)

Type: Subtask

Assignee: [Dominik Horvath](#)

Time Spent: Not Specified

Original estimate: Not Specified

Attachments: output.csv

Comments

Comment by [Dominik Horvath](#) [10/Oct/21]

Vytvoreny custom parser: <https://github.com/Dominik1799/gnewsParser>

Vsetky kniznice, ktore sme nasli, pracovali len ako wrapper nad RSS streamom. Vyuzivali funkcie, ktore su pre nas problem zbytocne, a boli preto velmi pomale.

Vlastny parser uskutocnil 1827 requestov a ziskal 35772 zaznamov za 93 sekund processing time a cca 20 minut realneho casu. Vyhľadavacia query ziskala vsetky clanky v casovom období od 10.10.2016 do 10.10.2021 na klucove slova "Jeffrey Epstein".

[AMS-20] [Analyza parametrov Google News](#) Created: 05/Oct/21 Resolved: 13/Oct/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Adam Šípka
Time Spent:	Not Specified
Original estimate:	Not Specified

Comments

Comment by [Adam Šípka](#) [13/Oct/21]
<http://books.gigatux.nl/mirror/googlehacks/0596008570/googlehks2-CHP-4-SECT-3.html>

[AMS-21] [Vytvorit Slack server](#) Created: 05/Oct/21 Resolved: 11/Oct/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Táňa Poláková
Time Spent:	Not Specified
Original estimate:	Not Specified

Description

Pre efektívnejšiu komunikáciu a strenutia v online priestore.

Comments

Comment by [Táňa Poláková](#) [11/Oct/21]

Vytvorila som Team na platforme Microsoft Teams

[AMS-3] [Zoznam trestnych cinov, ktore budu pouzite ako queries na google news](#) Created: 01/Oct/21 Resolved: 04/Oct/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Jakub Müller
Time Spent:	Not Specified
Original estimate:	Not Specified
Attachments:	 list_of_crimes.txt

Description

Zoznam zlocinov v anglictine, ktore budu pouzite ako query v Google News.
Schválené vedúcim.

[AMS-5] [Vysoka architektura](#) Created: 01/Oct/21 Resolved: 03/Oct/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	David Silady
Time Spent:	Not Specified
Original estimate:	Not Specified
Attachments:	 ams_high_arch.png

Description

Popripade k nodom v grafe aj na com to bezi (databaza, server, technologia, atd.)

Poznámka vedúceho: minimalizovať náklady na vývoj aby sa vyskúšala životnosť stránky; môže to všetko bežať na jednom serveri.

[AMS-6] [Specifikacia požiadaviek](#) Created: 01/Oct/21 Resolved: 03/Oct/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Jakub Hlavačka
Time Spent:	Not Specified
Original estimate:	Not Specified

Description

Poznámka vedúceho:

- Nefunkčné požiadavky:
 - o Fungovať by to malo aj bez prihlásenia
 - o Mohlo by to vyhodiť niečo aj bez prihlásenia
 - o Registrácia ľudí odrádza
 - o Pre neprihlásených vyhodit' len prvých 5 článkov
 - o Responzívny design – web aj mobil
 - o História, archív pre používateľov – venovať sa až neskôr
- Funkčné požiadavky:
 - o Pridať aj právnické osoby
 - o Pri scrapovaní musíme rátať aj s tým, že nejaký článok sa môže objaviť aj po našom scrapovaní – čo s vecami, ktoré sa mohli zmeniť

Comments

Comment by [Jakub Hlavačka](#) [03/Oct/21]

https://docs.google.com/document/d/18XLz0RXSFB50VAKdrydNxtY9zznVDNJK0wyjOS_NG0s/edit?usp=sharing

[AMS-7] [Prieskum ziskavania dat, praca s kniznicou](#) Created: 01/Oct/21 Resolved: 05/Oct/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Adam Šípka
Time Spent:	Not Specified
Original estimate:	Not Specified


Description

Kniznica Google News.

Poznámka vedúceho: mali by sme tu venovať veľa pozornosti, aby sme to celé nemuseli preprogramovať

[AMS-8] [Pripady pouzitia](#) Created: 01/Oct/21 Resolved: 05/Oct/21

Status:	Done
Project:	Adverse Media Screening

Type:	Story
Assignee:	Dominik Horvath
Time Spent:	Not Specified
Original estimate:	Not Specified
Attachments:	 Snímka obrazovky 2021-10-05 124128.png

Description

1. User opens AMS web
2. User enters search query - a name of a person
3. System returns list of news articles that contain this name
4. user clicks on one of the results
5. system redirects the user to the news source

Poznámka vedúceho: je to ok.