

Automatizovaná analýza a archivace dat z webu

Tomáš Kocman*



Abstrakt

Archivace dat z webu je užitečná pro ty případy, ve kterých si chceme udržovat měnící se informace o nějakém subjektu v čase. Tato práce umožňuje automatizovat archivaci webových stránek, ovšem jen těch, které splňují určitá pravidla – data na nich obsažená vyhovují definovaným regulárním výrazům. Výsledkem práce je platforma, kterou lze konfigurovat takovým způsobem, aby prohledávala a archivovala webové stránky podle různých strategií. Mějme například instituci jako muzeum nebo knihovnu, která by chtěla ukládat historii určitých dokumentů na webu. S platformou lze jednoduše automatizovaně navštívit všechny stránky na daném webu a pokud tyto stránky splňují definovaná pravidla, platforma provede jejich zálohu. V oblasti kyberkriminality například vyšetřovatelé znají webové stránky, popřípadě fórum, kde pachatel prováděl trestnou činnost. Potom mohou platformu využít k nalezení důkazního materiálu – internetovou přezdívku pachatele, obsah zpráv a další pro soud cenné informace.

Klíčová slova: Kyberkriminalita — Archivace webu — Dolování dat

Příložené materiály: [Archivované stránky](#)

*xkocma04@stud.fit.vutbr.cz, *Fakulta informačních technologií, Vysoké učení technické v Brně*

1. Úvod

Během posledních několika let se *World Wide Web* stal doslova napumpovaný informacemi a v některých případech se stalo velmi obtížným najít konkrétní informaci. Jakožto řešení tohoto problému vznikla disciplína *web mining* [1].

Pod pojmem *web mining* si lze představit aplikaci různých technik pro dolování dat. Je to proces objevování potenciálně užitečných a doposud neznámých informací. *Web mining* se využívá pro:

- získávání relevantních informací,
- vytváření nových znalostí z relevantních dat,
- personalizaci informací,

- nauku o zákaznících a obecně o uživateli, pokud jde o nějaký typ služby.

Tato práce přináší platformu, jejíž úlohou je automatizovaně prohledávat webové stránky a jestliže daná stránka spadá do našeho okruhu zájmu, provede její archivaci. Platforma je plně konfigurovatelná co se týče strategií průchodu webovými stránkami a je plně modulární.

Nejprve budou v kapitole 2 uvedeny informace z oblasti dolování obsahu z webu, proč je to zajímavé a hlavní body důležité k orientaci v této doméně. Následující kapitola 3 obsáhne specifika vyšetřování trestné činnosti v kyberprostoru spolu s definicí digitální stopy, která je pro tuto oblast esenciální. V kapitole 4 už bude

14
15
16
17
18
19
20
21
22
23
24
25
26
27

28 detailněji popsán koncept celé platformy s činnostmi
29 každé jednotlivé části, kde dominujícími tématy budou
30 Scrapy a Lemmiwinks. Poslední kapitola 5 uvede
31 různé způsoby, jak lze platformu nasadit a krátký výčet
32 vhodných použití.

33 2. Sběr dat na internetu

34 Dolování webového obsahu je podmnožina tohoto po-
35 jmu. Se zvyšující se složitostí webů se rovněž enormně
36 zvyšuje množství informací na nich obsažených. Tudíž
37 extrakci dat, kterou by uživatel mohl požadovat se stala
38 náročným a často zdoluhavým úkolem [1]. Výsledkem
39 je to, že se dolování dat z webů stalo základní tech-
40 nikou pro extrakci hodnotných informací z internetu.
41 *Web mining* [2] je dále rozčleněn do tří kategorií –
42 dolování webového obsahu, dolování struktury webu
43 a dolování užívání webu:

- 44 • při dolování webového obsahu zkoumáme ob-
45 jekty jako text, obrázky nebo multimedia,
- 46 • při dolování webové struktury pracujeme na
47 základě struktury webu, kdy například hledáme
48 odkazy specifikované pomocí URL,
- 49 • v případě dolování webového užívání jsou obla-
50 stí zájmu logovací soubory, které obsahují na-
51 vigační vzorce uživatelů, čili jak se uživatelé
52 pohybovali na daném webu.

53 2.1 Dolování obsahu z webu

54 Dolování obsahu z webu má svou vlastní taxonomii.
55 Pro tuto práci je podstatné dolování webového obsahu,
56 konkrétně textu obsaženého na specifických stránkách.
57 Obecně jde o běžnou techniku, kdy prohledáváme web
58 skrze jeho obsah. Rovněž vyhledávací procesory dělají
59 navíc ke své činnosti dolování webového obsahu.

60 V oblasti dolování webového obsahu jsou pro tuto
61 práci stěžejní převážně techniky dolování nestruk-
62 turovaného obsahu (převážně textu). Dolováním nestruk-
63 turovaných dat získáváme neznámé informace. Dolo-
64 vání textu spočívá v extrakci dříve neznámých infor-
65 mací z různých textových zdrojů [2].

66 V práci se užívá techniky extrakce informací, která
67 používá vzory k získání shody textu. Vhodnými meto-
68 dami jsou vyhledávání ve slovníku (hledání klíčových
69 slov) a regulární výrazy (vyhledání celých frází dle
70 specifického předpisu). Tato technika je vhodná při
71 velkých objemech dat [2].

72 Když se vytváří nástroj na sběr webových dat,
73 jsou zde citlivé oblasti, na které musí vývojář brát
74 zřetel. Nezodpovědné sbírání dat může být pro druhou
75 stranu minimálně otravné, ovšem může to v některých
76 případech hraničit až s ilegality. Je potřeba se vy-

varovat dvěma primárními věcmi – útoky DoS (odepře- 77
ní služby) a porušení autorských práv [3]. 78

3. Vyšetřování kyberkriminality 79

80 Za kyberkriminalitu lze považovat jakékoli protiprávní 80
jednání. Kyberterorismus úzce souvisí s tímto pojmem, 81
jelikož značí takové teroristické aktivity v kyberpros- 82
toru, které narušují počítačové sítě a zařízení. Při 83
útocích tohoto typu může docházet k lidským úmrtím 84
nebo k závažným ekonomickým ztrátám, jejichž důsle- 85
dky jsou jen těžko předvídatelné. 86

87 Podle bezpečnostního auditu ministerstva vnitra 87
ČR [4] zde hrají významnou roli tzv. nová média 88
(média založená na digitálním kódování dat – soft- 89
ware i hardware). Kyberkriminalitu nelze považovat 90
za pouhý hypotetický fenomén. Útoky probíhají v čím 91
dál větším měřítku nejen na běžné uživatele internetu, 92
čili jednotlivce, ale rovněž na celé státy (například 93
pro manipulaci s politickou situací). V současnosti 94
však velkou část útoků a incidentů, často mediálně 95
popisovaných a prezentovaných jako kyberkriminalita, 96
můžeme označit spíše za využívání kyberprostoru (in- 97
ternetu) teroristy. Teroristické organizace zatím ne- 98
jspíše nedisponují kapacitami a schopnostmi k usku- 99
tečnění kybernetických útoků s vážnými dopady [4]. 100
Ovšem není obtížné tyto schopnosti a kapacity nakou- 101
pit ve formě služby. Například Islámský stát dokázal 102
provést kybernetické útoky (avšak nikterak sofistiko- 103
vané), které jiné teroristické organizace nebyly schopny 104
dlouho uskutečnit. 105

3.1 Digitální stopa 106

107 Stále častěji se součástí důkazního materiálu, a to ne- 107
jen v oblasti počítačové kriminality, stávají rovněž 108
digitální stopy [5]. Definicí digitální stopy je mnoho, 109
různí autoři často používají synonyma pro označení 110
digitální stopy, nám zde však postačí krátká výstižná 111
definice – *digitální stopa je informace zanechaná od 112
uživatele v prostředí internetu nebo jako součást sou- 113
borů*. Pojem digitální stopa je pro účely této práce 114
stěžejní, jelikož se snaží právě o sběr aktivních digitál- 115
ních stop. Aktivní informace, které po sobě uživatelé 116
na internetu zanechávají jsou: 117

- 118 • profily nebo příspěvky zanechané na sociálních 118
sítích, 119
- 120 • e-maily, sms zprávy, historie chatu, 120
- 121 • různé typy úředních údajů. 121

122 Jak je jistě patrné, mezi aktivní se řadí ty informace, 122
které o sobě uživatel zveřejní prostřednictvím různých 123
služeb vědomě a dobrovolně [5]. Naopak pasivní infor- 124
mace vznikají v prostředí internetu bez našeho přímého 125

126 záměru. Realita je taková, že jakákoliv aktivita v on-
127 line prostředí může být zaznamenána a uložena. Tyto
128 informace se často zneužívají například pro:

- 129 • krádež osobních informací (údaje z kreditních
130 karet, rodné číslo, e-mailová adresa),
- 131 • kyberšikana,
- 132 • kyberstalking,
- 133 • zdroj informací pro personalisty,
- 134 • sledování návyků uživatelů (realizováno třetími
135 stranami, sběrateli dat a reklamními společno-
136 stmi).

137 Existují však způsoby a nástroje pro správu a ochranu
138 digitálních stop. Pro správu aktivních stop je účinné:

- 139 • používat více přihlašovacích jmen,
- 140 • rozvážně (nejlepe vůbec) publikovat fotografie,
141 videa a osobní údaje,
- 142 • vhodně nastavit soukromí, obzvláště u sociál-
143 ních sítí, které poskytují základní konfiguraci,
- 144 • vhodně nastavit zabezpečení svého internetové-
145 ho prohlížeče,
- 146 • Me on the web¹ (služba, která sleduje nově zve-
147 řejněné informace o uživateli v určité oblasti).

148 4. Platforma pro sběr důkazů

149 Ačkoli je v rámci práce vytvořena další, odlehčená,
150 verze platformy podporující systémy Windows, popis
151 platformy se bude opírat o obrázek 1, čili o plnohod-
152 notnou Unixovou verzi. Jak lze vidět na obrázku, plno-
153 hodnotná verze se skládá z pěti částí:

- 154 • Scrapy je knihovna pro sběr a analýzu webových
155 dat. Umožňuje dynamickou konfiguraci pavouků
156 a tzv. potrubí (viz. sekce 4.1), přes která tečou
157 stažená data a v každém tomto bodě lze s daty
158 provádět odlišné operace. Procesů Scrapy může
159 být v rámci platformy více, podle aktuálního
160 zadání.
- 161 • Redis je distribuovaná fronta pracující pouze
162 v rámci paměti nezávisle od ostatních kompo-
163 nent platformy. Zde stačí pouze jeden databá-
164 zový datový typ list, který slouží jako fronta.
165 Procesy Scrapy vkládají do fronty data a pro-
166 cesy Lemmiwinks z fronty tato data odebírají.
167 Do fronty se ukládají celé HTML dokumenty
168 spolu s metadaty.
- 169 • Lemmiwinks je webový archivátor, ukládající
170 výstupní archivy ve formátu MAFF. Proces z Re-
171 disu načítá celé HTML dokumenty a ty rekurzi-
172 vně archivuje. V rámci platformy může být in-
173 stancí procesu Lemmiwinks více, jelikož fronta

v Redisu se může plnit mnohem rychleji, nežli
Lemmiwinks dokáže data konzumovat.

- 174 • PostgreSQL je perzistentní databáze, která přiřa-
175 zuje k jednotlivým doménám archivovaných we-
176 bových stránek cesty k jejich uloženému MAFF
177 archivu. Schéma databáze je obohaceno o meta-
178 data související s archivovanou webovou strán-
179 kou.
180 • Management slouží jako aplikační programové
181 rozhraní celé platformy. Skrze tohle REST API
182 je možné manipulovat s jednotlivými částmi
183 platformy – spouštět/ukončovat procesy Scrapy
184 nebo Lemmiwinks a získávat lokaci archivu dle
185 zadaného dotazu. V rámci obrázku 1 je Manage-
186 ment část spojena s každou z ostatních částí.
187
188

Nejdůležitějšími částmi celé platformy jsou Scrapy pro
sběr a analýzu webových dat a Lemmiwinks jakožto
webový archivátor.

4.1 Scrapy 192

Scrapy je robustní webová Python knihovna pro dolo-
vání dat z různých zdrojů [6]. Z vysokoúrovňového
pohledu Scrapy exceluje obzvláště při dvou případech
užití:

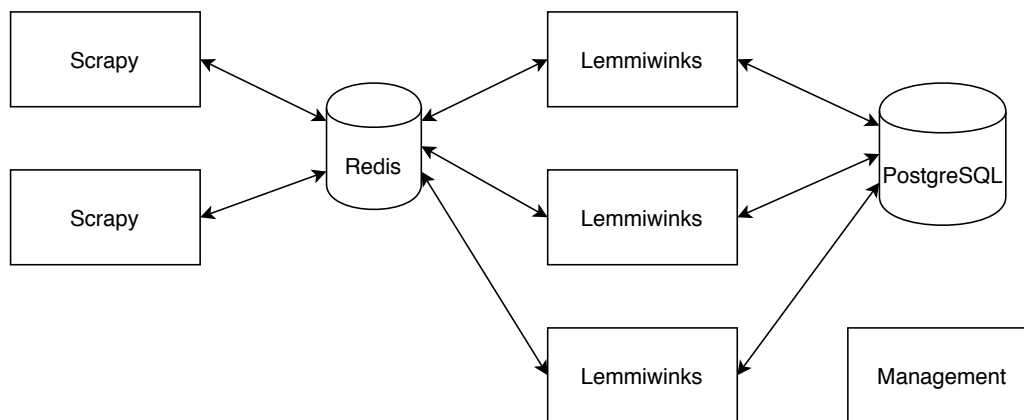
- 193 • Když běžný uživatel webu požaduje stáhnout
194 určitá data ze stránky, kterou zrovna prohlíží a
195 libovolně data formátovat. Například je uložit ve
196 formátu JSON nebo CSV či je uložit do databáze
za účelem offline prohlížení nebo provedení dal-
ších výpočtů.
- 201 • Pokud si uživatel přeje kombinovat data z růz-
202 ných zdrojů a extrahovat je.

S knihovnou Scrapy jsme schopni při jedné konfiguraci
provést úlohy, na které bychom s jinými nástroji nebo
knihovnamy potřebovali mnoho tříd, rozšíření a kon-
figurací [6]. Z pohledu programátora stojí za zmínku
event-based architektura. Ta umožňuje vytvářet kaská-
dy operací, které mohou čistit, formátovat, obohacovat
nebo data ukládat například do databáze, přičemž neza-
znamenané žádnou degradaci výkonu.

Scrapy je již hotový open-source projekt, který
byl do platformy převzat. Následně jej bylo potřeba
doplnit o definici konkrétních pavouků sbírající data,
potrubí, která data modifikují, popřípadě ukládají do
databáze a vhodné konfigurační parametry, které činí
sběr jednoduchým a efektivním.

Základními stavebními bloky aplikací využívajících
knihovnu Scrapy jsou pavouci a potrubí. Pavouci
vytváří HTTP dotazy, zpracovávají HTTP odpovědi
a poté generují jednotlivé prvky a další dotazy. Každý

¹<https://myaccount.google.com/dashboard>



Obrázek 1. Architektura plnohodnotné platformy pro systémy Unix.

223 z prvků, které jsou pavoukem vygenerovány jsou ná- 261
 224 sledně zpracovány sekvencí potrubí (potrubí si lze 262
 225 v tomto kontextu představit jako posloupnost funkcí 263
 226 nebo operací, kdy každá z těchto funkcí přijímá prvek 264
 227 jakožto parametr). Typicky tyto funkce nějakým způ- 265
 228 sobem modifikují vstupní prvek a posílají jej dále 266
 229 k další funkci v potrubí jednoduše tím, že je prvek 267
 230 použit jako návratová hodnota funkce. Příležitostně 268
 231 (například při detekci duplikátního prvku) je požadova- 269
 232 né chování zahození prvku. V takovém případě další 270
 233 části potrubí tento prvek nedostanou a tudíž je zas- 271
 234 taveno následující zpracování. 272

235 Při sběru dat je podstatný HTML dokument, který 273
 236 Scrapy uloží do Redis databáze spolu s metadaty, které 274
 237 budou v konečném důsledku zpropagovány až do Post- 275
 238 greSQL. Lemmiwinks bude následně archivovat we- 276
 239 bovou stránku na základě tohoto HTML dokumentu, 277
 240 který získal Scrapy a to kvůli zachování integrity dat. 278
 241 Webové stránky se totiž mohou měnit každou sekundu 279
 242 a pokud by Scrapy předal archivačnímu rámci pouze 280
 243 URL adresu, kterou má archivovat, mohla by být in- 281
 244 tegrity dat porušena. Metadata, která Scrapy ukládá 282
 245 jsou: 283

- 246 • URL adresa,
- 247 • název běžící úlohy,
- 248 • název pavouka,
- 249 • jméno serveru, na kterém úloha běží (platforma 284
 250 může fungovat v distribuovaném prostředí),
- 251 • datum a čas,
- 252 • regulární výraz použitý pro běžící úlohu,
- 253 • seznam řetězců, které vyhovují regulárnímu vý- 285
 254 razu, oddělených středníkem,
- 255 • HTTP hlavička,
- 256 • HTML dokument.

257 4.2 Lemmiwinks

258 Lemmiwinks je programový rámec poskytující funkci- 286
 259 onality získávání dat z webů a jejich archivaci [7]. 287
 260 Velkým přínosem tohoto rámce je možnost zpracovat 288
 289

261 také webové stránky s dynamickým obsahem. Lemmi- 262
 263 winks archivuje celé webové stránky a to do formátu 264
 265 MAFF (Mozilla Archive Format), jehož předností je 266
 267 schopnost archivovat více oteřených záložek do jed- 268
 269 noho výsledného dokumentu. 270

271 Lemmiwinks je již hotový open-source projekt, 272
 273 který byl do platformy převzat. Následně jej bylo 274
 275 potřeba doplnit o nové třídy, které využívají archivační 276
 277 funkcionality, třídy operující s Redisem a PostgreSQL 278
 279 a v neposlední řadě modifikovat způsob provádění 280
 281 celého rámce tak, aby dokázal pracovat s asynchronní 282
 283 smyčkou událostí. 284

285 Architektura Lemmiwinks je modulární, umožňuje 286
 287 tudíž přidávat libovolně nové moduly, případně měnit 288
 289 implementaci těch stávajících [7]. Jelikož je architek- 286
 287 tura modulární, má generický design, který definuje 288
 289 a používá rozhraní pro řešení částečných problémů. 289
 Jako programovací jazyk je použit Python. Určité části 286
 rámce jsou paralelní kvůli zvýšení výkonu a lepšímu 287
 využití zdrojů. Tento paralelismus není implemen- 288
 tován pomocí vláken, ale po vzoru Scrapy byla využita 289
 asynchronní architektura, kterou zajišťuje standardní 286
 knihovna asyncio (Scrapy ovšem využívá asynchronní 287
 knihovnu Twisted, která je méně výkonná a aplikačně 288
 náročnější než asyncio). 289

286 5. Možnosti nasazení

287 Způsobů a míst, kde nasadit takovou platformu je ne- 288
 289 spočet, proto bude uvedena jen jejich malá, pro čtenáře 289
 zajímavá, podmnožina: 286

- 290 • vyhledání nelegálních torrentů podle info hashů,
- 291 • nalezení čísel odzicených bankovních karet,
- 292 • hledání adres sítě Tor,
- 293 • sběr informací o uživateli (jeho celková aktivita 294
 295 na webu),
- 296 • hledání prodeje nelegálního zboží,
- 297 • vyhledání telefonních čísel,

- 297 ● archivace dokumentů za účelem budování histo- 348
- 298 rie webu, 349
- 299 ● archivace cen určitého zboží v čase.

300 Platforma je aktuálně dokončena a probíhá její testo- 351

301 vání na úrovni výkonu a potřeby perzistentního pamě- 352

302 ťového uložení. Regresní testování probíhá na strojově 353

303 generovaných webových stránkách s předdefinovaným 354

304 obsahem. Nejprve je nakonfigurován Scrapy, aby vy- 355

305 bral podle regulárního výrazu pouze specifické stránky, 356

306 které uloží do databáze Redis. Ty následně Lemmi- 357

307 winks vyčítá a archivuje je. Poté je automatizovaně 358

308 porovnává originální stránka z webového serveru se 359

309 stránkou uloženou v archivu.

310 Dalším typem testování jsou jednotkové testy, které 361

311 jsou specifické pro Scrapy i Lemmiwinks. Scrapy má 362

312 svůj vlastní systém testování, kdy každému pavouku 363

313 jsou definovány kontrakty. Kontrakt deklarativně popi- 364

314 suje, která webová stránka se má stáhnout a jaké ob- 365

315 jekty na ní má pavouk najít, aby byl kontrakt splněn. 366

316 Lemmiwinks je poté testován pomocí Python knihovny 367

317 *pytest*.

318 Přínos představené práce spočívá v robustnosti 368

319 řešení, které je plně konfigurovatelné, modulární a je 369

320 od začátku tvořeno s představou, že bude fungovat dis- 370

321 tribuovaně. V porovnání s existujícími nástroji jako je 371

322 například wget nebo curl se práce liší převážně v tom, 372

323 že obsahuje hotový webový sběrač dat, který automa- 373

324 tizovaně navštívuje jednotlivé stránky daného webu 374

325 a na základě splnění předem definovaných podmínek 375

326 stránku archivuje. V porovnání se zmíněnými nástroji 376

327 se diametrálně liší i samotná archivace. Zatímco curl 377

328 nebo wget pouze stáhnou HTML dokument, platforma 378

329 se postará o kompletní archivaci se všemi zdroji, které 379

330 stránka obsahuje (obrázky, CSS styly a další doku- 380

331 menty), tudíž při opětovném otevření archivu vypadá 381

332 stránka totožně jako v čase archivace. Další výhodou, 382

333 která je v dnešní době spíše samozřejmostí je schop- 383

334 nost archivace webových stránek s dynamickým we- 384

335 bovým obsahem.

336 6. Závěr

337 Tato práce popisuje řešení pro automatizované stažení, 385

338 analýzu a archivaci webových stránek. Výsledkem 386

339 práce je platforma určená pro vyšetřovatele a bezpečno- 387

340 stní experty České republiky. Platforma se skládá ze 388

341 čtyř částí – Scrapy pro stažení a analýzu dat, databáze 389

342 Redis sloužící jako fronta, Lemmiwinks pracující jako 390

343 archivační nástroj webových stránek a databáze Post- 391

344 greSQL pro perzistentní uložení dat. Všechny kompo- 392

345 nenty pracují v kontejnerech Docker. Cílovou skupinou 393

346 této práce jsou bezpečnostní experti, ale i běžní uživatelé, 394

347 kteří potřebují z nějakého důvodu vyhledat na webu

požadovanou informaci a zaarchivovat stránku, na 348

349 které byla informace nalezena.

Práce byla zahrnuta do projektu TARZAN. Jde 350

351 o integrovanou platformu pro zpracování digitálních 352

353 dat z bezpečnostních incidentů. Projekt je zaměřen na 354

355 analýzu a detekci nových forem kybernetické krimina- 356

357 lity především v prostředí mobilních a komunikačních 358

359 aplikací a v prostředí internetu věcí. Cílem projektu 360

361 je výzkum nových technologií a metod pro efektivní 362

363 vyšetřování bezpečnostních incidentů. Výsledky bu- 364

365 dou demonstrovány na případech z praxe, například 366

367 na detekci provozu P2P sítí, bezpečnostní analýze mo- 368

369 bilních zařízení či řešení incidentů v oblasti Bitcoinů. 369

370

Poděkování 361

Rád bych poděkoval svému vedoucímu panu Ing. Li- 362

363 boru Polčákovi Ph.D. a panu Ing. Viliamu Serečunovi 364

365 za odborné rady a čas, které mi při tvorbě práce poskytli. 366

Literatura 365

- [1] Faustina Johnson. *Web Content Mining Tech- 366*
- niques: A Survey. International Journal of Com- 367*
- puter Applications*, 47(11), 7 2012. 368
- [2] Bing Liu. *Web Data Mining: Exploring Hyper- 369*
- links, Contents, and Usage Data*. Springer Science 370
- & Business Media*, 2011. 371
- [3] Katharine Jarmul, Richard Lawson. *Python Web 372*
- Scraping*. Packt Publishing, 2017. 373
- [4] Ministerstvo vnitra ČR, odbor bezpečnostní poli- 374
- tiky a prevence kriminality. Audit národní 375*
- bezpečnosti*. [Online; navštíveno 27.09.2018]. 376
- [5] Jan Kolouch. *CyberCrime*. CZ.NIC, 2016. 377
- [6] Dimitrios Kouzis-Loukas. *Learning Scrapy*. Packt 378
- Publishing*, 2016. 379
- [7] Viliam Serečun. *Automatizovaná rekonstrukce 380*
- webových stránek*. Master's thesis, Vysoké učení 381
- technické v Brně, Fakulta informačních tech- 382*
- nologií*, 2018. 383