

Využití moderních technologií pro práci s obsahem webového archivu

Jan Lehečka, 21.10. 2019

Workshop výzkumného projektu
“Vývoj centralizovaného rozhraní pro
vytěžování velkých dat z webových archivů”

KATEDRA
KYBERNETIKY



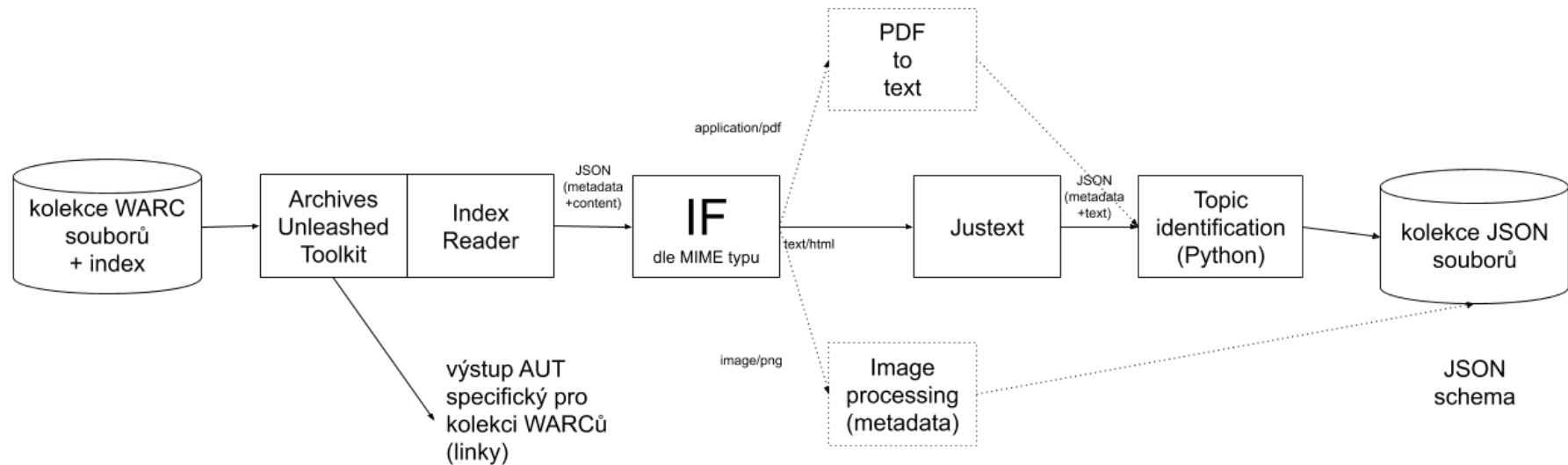
FAKULTA APLIKOVANÝCH VĚD
ZÁPADOČESKÉ UNIVERZITY
V PLZNI



Osnova prezentace

- Schéma řešení a použité technologie
- Specifikace intermediary formátu
- Ukázka extrakce textu z webové stránky
- Ukázka identifikace tématu webové stránky

Schéma řešení a použité technologie



[Webarchiv](#)

Archives Unleashed

PySpark

python

scikit learn

APACHE HBASE

hadoop + APACHE Spark

Specifikace intermediary formátu

- *intermediary formát* = formát jednotlivých záznamů během zpracovávání webového archivu
- definuje položky pro surová stažená data a metadata
- specifikace definována pomocí JSON schéma
- každá webová stránka = 1 JSON objekt, do kterého postupně jednotlivé algoritmy ukládají získaná metadata

The screenshot shows a JSON schema viewer interface. At the top, there are tabs for 'JSON', 'Surová data', and 'Hlavičky'. Below the tabs are buttons for 'Uložit', 'Kopírovat', 'Sbalit vše', 'Rozbalit vše', and a search filter 'Filtr JSON'. The main content area displays a JSON schema for 'schema.json#'. The schema is a JSON object with the following properties:

id:	"schema.json#"
\$schema:	"http://json-schema.org/draft-07/schema#"
type:	"object"
title:	"Intermediary format"
description:	"This is a JSON schema of...ormat for WebArchiv NK."
required:	[-]
properties:	
id:	[-]
content:	[-]
plain-text:	[-]
plain-text-tokens:	[-]
plain-text-sentences:	[-]
urlkey:	[-]
timestamp:	[-]
url:	[-]
mime-type:	[-]
response-code:	[-]
digest:	[-]
redirect-url:	[-]
robot-meta-tags:	[-]
warc-offset:	[-]
warc-record-size:	[-]
warc-filename:	[-]
rec-headers:	[-]
http-headers:	[-]
title:	[-]
headlines:	[-]
links:	[-]
language:	[-]
web-page-type:	[-]
topics:	[-]
sentiment:	[-]

Ukázka extrakce textu z webové stránky



jusText

Neděle 20. října 2019, svátek má Vendelín

Vyzkoušet Premium zdarma

Premium

Přihlásit Můj účet

iDNES.cz / Sport

V semifinále MS se střetnou ragbisté Nového Zélandu a Anglie

Ostatní sporty

Dalších 8 fotografií v galerii

Australský ragbista Samu Kerevi (u míče) v duelu s Anglii. | foto: AP

říj 19 2019

Ragbisté Nového Zélandu deklasovali na mistrovství světa v Japonsku ve čtvrtfinále 46:14 Irsko a udělali další krok k třetímu titulu za sebou. Soupeřem All Blacks v semifinále bude příští sobotu Anglie, jež v sobotu porazila podobně jednoznačně 40:16 Austrálii.

Vstoupit do diskuse

Irsko sice z předchozích tří zápasů Nový Zéland dvakrát porazilo, v Tokiu ale vládli trojnásobní mistři světa. Položili sedm pětěk, z toho první dvě Aaron Smith. Irové, kteří byli před turnajem světovými jedničkami a usilovali o první semifinále na MS v historii, vymazali nulu na kontě až v závěrečné desetimínutovce za rozhodnutého stavu.

Anglie se v souboji bývalých mistrů světa v Oitě dostala do vedení 14:3 v prvním poločase díky dvěma pětkami Johnnyho Maye během tří minut. Australané se sice krátce po pauze zásluhou pětky Mariky Korobeteho přiblížili na rozdíl jedného bodu, na to ale okamžitě zareagoval pětkou Kyle Sinckler a Anglie už zbytku utkání dominovala.

Australané nepostoupili do semifinále teprve poťetí v devíti ročnících MS a ve všech případech podlehl Anglii. Ta přešla přes čtvrtfinále poprvé od finále v roce 2007.

V neděli určí druhou semifinálovou dvojici zápasy Wales - Francie a Japonsko - Jihoafrická republika.

Mistrovství světa v ragby v Japonsku

Čtvrtfinále

Anglie - Austrálie 40:16 Nový Zéland - Irsko 46:14

Autoři: ČTK, iDNES.cz

Související

Nic než zlato. Novozélandané vědí, že jejich země jinou medaili neuznává

Japonci vyzývají další zázrak. A co Irové? Fascinující podivaná začíná

Sport v roce 2019

20. 9. - 2. 11. MS v ragby (Japonsko)

26. 12. - 5. 1. MS hokejistů do 20 let (Ostrava a Třinec)

Témata: MS v ragby 2019 v Japonsku, Ragby

Vstoupit do diskuse

Sdílet na Facebooku

demo: <https://nlp.fi.muni.cz/projects/justext/>

Ukázka identifikace tématu webové stránky

Kurdové a Turecko se vzájemně obvinili z nedodržování příměří v Sýrii



KURDŮ VYHOZÍVÁJÍCÍ Z MĚSTA NA HRANICÍ TURECKA SE SÝRIÍ. TURECKO ZAČALO 9. ŘÍJNA. S. | FOTO: AP PHOTO

DAMASEK/ANKARA Turecko a syrští Kurdové se navzájem obvinili z nedodržování příměří na severu Sýrie. Bylo ujednáno ve čtvrtek a má platit 120 hodin. Kurdové v té době mají stáhnout své oddíly z území u turecké hranice. Turecký prezident Recep Tayyip Erdogan v pátek pohrozil, že turecká armáda ofenzívu v úterý večer obnoví, jestliže Kurdové zónu nevyklidí. Turecký ministr obrany v sobotu tuto hrozbu zopakoval.

V sobotu turecké ministerstvo obrany oznámilo, že turecká armáda dohodu vyjednanou ve čtvrtek v celém rozsahu plní. „Teroristé přesto v minulých 36 hodinách zanořili čtrnáctkrát,“ uvádí se v prohlášení úřadu, který slovem teroristé označuje syrské kurdské oddíly YPG.

Ministr obrany Hulusi Akar řekl, že Turecko přenášelo boje, aby se mohli kurdští bojovníci stáhnout. „Jejich zbraně budou zabaveny a jejich posty zničeny. Pokud se to nestane, budeme dále bojovat,“ řekl Akar.

Stanice Al-Džazíra má informace, že do vyjednávání s Kurdy vstoupili syrští povstalci, které Turecko podporuje od začátku syrské války a kteří se turecké ofenzívy zúčastnili. Podle Al-Džazíry jsou nyní tyto bojovníci u Rás al-Ajnu a snaží se přimět arabsko-

Neslovník Claudiu Johnson douhá, že trasování příměří nezastaví obranu. Sílněji ale komplikuje pozemňovací návrh



ADVENTNÍ TRHY A PRŮVODY ČERTŮ

Ukázat zájezdy

Ceské kormidlo

novoc.cekorkmidlo.cz

jusText

- Plain text -

Damašek/Ankara Turecko a syrští Kurdové se navzájem obvinili z nedodržování příměří na severu Sýrie. Bylo ujednáno ve čtvrtek a má platit 120 hodin. Kurdové v té době mají stáhnout své oddíly z území u turecké hranice. Turecký prezident Recep Tayyip Erdogan v pátek pohrozil, že turecká armáda ofenzívu v úterý večer obnoví, jestliže Kurdové zónu nevyklidí. Turecký ministr obrany v sobotu tuto hrozbu zopakoval.

V sobotu turecké ministerstvo obrany oznámilo, že turecká armáda dohodu vyjednanou ve čtvrtek v celém rozsahu plní. „Teroristé přesto v minulých 36 hodinách zaútočili čtrnáctkrát,“ uvádí se v prohlášení úřadu, který slovem teroristé označuje syrské kurdské oddíly YPG.

Ministr obrany Hulusi Akar řekl, že

Get labels



- Topics -

#	label	score
1	turecko	1.000
2	sýrie	0.970
3	nepokoje, konflikty a válka	0.881
4	diplomacie	0.426
5	zločin, zákon a spravedlnost	0.107
6	národní vláda	0.053
7	irák	0.028
8	ozbrojené síly	0.022
9	eu	0.020
10	místní úřad	0.012



Děkuji za pozornost