



# Part 1: Latest developments in DBpedia Part 2: Entityclassfier.eu NER

#### Milan Dojčinovski

**KEG Seminar 5 December, 2019** 



Initially started in 2007 in Leipzig, Germany

2008 - DBpedia is a crowd-sourced **community effort** to extract structured **information from Wikipedia** and make this information **available on the Web**.

2018 - New mission: Global and unified access to knowledge graphs

- Original definition still holds true
- Platform (matchmaking) to integrate your data with all other data
- DataBus https://databus.dbpedia.org

# **DBpedia allows to access data**





Connecting data is about connecting people and organisations

#### **Czech DBpedia**

Po směru hodinových ručiček počínaje obrázkem nahoře: Pražský hrad, výškové budovy na Pankráci, Malá Strana, Staroměstské náměstí, Karlův most, Národní divadlo



Znak Prahy Praga Caput Rei heslo: publicae<sup>[p. 1]</sup> (dříve Praha matka měst) status: hlavní město, zároveň krai a statutární město historická země: Čechy LAU 2: CZ0100 554782 kraj (NUTS 3): Hlavní město Praha (07010)

	(02010)
okres (LAU 1):	Hlavní město Praha (CZ0100)
ISO 3166-2:CZ:	CZ-PR
Státní poznávací značka:	A
poštovní směrovací číslo	100 00-199 00
katastrální výměra:	496 km²
obyvatel:	1 294 513 <sup>[1]</sup>
rozpočtové výdaje:	60 991 mil. Kč (2010) <sup>[2]</sup>
hustota zalidnění:	2581,7 obyvatel/km <sup>2</sup>
zemēpisná šířka:	50° 05' s. š.
zemēpisná délka:	14° 25' v. d.
nadmořská výška:	177–399 m n. m.
nejvyšší bod:	vrch Teleček mezi Sobínem a Chrášťany (399 m n. m.)
nejnižší bod:	hladina Vltavy u Suchdola (177 m n. m.)
počet městských obvodů	: 10
počet městských (správních) obvodů:	22

57

146

počet městských částí:

počet místních částí:





CATEGORIES

TYPES

GALLERY External Links

Born Here

Q



@ http://cs.dbpedia.org Cesky, čeština

#### Praha

Praha je hlavní a současně největší město Česka a 15. největší město Evropské unie. Leží mírně na sever od středu Čech na řece Vltavě, uvnitř Středočeského kraje, jehož je správním centrem, ale jako samostatný kraj není jeho součástí. Je sídlem velké části státních institucí a množství dalších organizací a firem. Sídlí zde prezident republiky, parlament, vláda, ústřední státní orgány a jeden ze dvou vrchních soudů.

cs.wikipedia.org/wiki/Praha



Property:	Value:	
prop-cs:aprHi°c :	13.4 (xsd:double)	
prop-cs:aprLo°c :	3.5 (xsd:double)	
prop-cs:aprPrecipMm :	38.2 (xsd:double)	
prop-cs:augHi°c :	23.5 (xsd:double)	
prop-cs:augLo°c :	13 (xsd:integer)	
prop-cs:augPrecipMm :	69.6 (xsd:double)	
prop-cs:další :	nafotografovali, sestavili, úvodem a rejstříkem opatřili Barbora a Marek Lašťovkovi @cs	-
prop-cs:decHi*c :	2.1 (xsd:double)	



# **Core Dataset Groups**



Available extractions 13 billion facts total (200 GB)

- Generic (automatic)
- Mappings-based (rule-based)
- Text
- Wikidata

Based on the Wikimedia XML dumps

# **Generic extraction**



- 132 languages, 30 datasets
- http://dbpedia.org/property/ properties
- https://en.wikipedia.org/w/index.php?title=Prague&action=edit
- <u>http://dbpedia.org/page/Prague</u>
- <u>https://github.com/dbpedia/extraction-framework/tree/master/core/src/</u> <u>main/scala/org/dbpedia/extraction/mappings</u>

# **Mappings based extraction**

- 40 languages, 6 datasets
- http://dbpedia.org/ontology properties

coordinates
subdivision\_type
subdivision\_name
established\_title
area\_urban\_km2

= {{coord|50|05|N|14|25|E|region:CZ|display=inline,title}}

= Country

- = [[Czech Republic]]
- = Founded
- = 298

Property Mapping (help)				
template property	area_urban_km2			
ontology property	areaUrban			
unit	squareKilometre			

http://dbpedia.org/resource/Prague

Geocoordinates Mapping (help)							
coordinates template property	coordinates						
Geocoordinates Mapping (help)							
longitude degrees template property	longd						
longitude minutes template property	longm						
longitude seconds template property	longs						
longitude direction template property	longEW						
latitude degrees template property	/ latd						
latitude minutes template property	latm						
latitude seconds template property	/ lats						
latitude direction template property	/ latNS						



### **Text extraction**



- 132 languages, 8 datasets
- Short and long abstracts
- Training data for text mining
- Fact extraction

Currently offline due to maintenance (refactoring)

# Wikidata extraction

- Same approach as for Wikipedia:
  - Generic and Mappings-based
- Mappings in JSON

- Allows unified access over Wikipedia and Wikidata
- + Wikidata has no ontology,DBpedia has 8 (DBO, Yago, Umbel,..)

+ Generic still extracts 584 million facts <sup>1</sup>



```
"P279": [
        "rdfs:subClassOf": "$getDBpediaClass"
],
"P625": [
        "rdf:type": "http://www.w3.org/2003/01/geo/wgs84_pos#Spatial1
    },
        "geo:lat": "$getLatitude"
    },
    {
        "geo:long": "$getLongitude"
    },
        "georss:point": "$getGeoRss"
```

# What's New in DBpedia



- 2+ years innovation phase (ongoing)
- Reengineered the released process -> DataBus
- More frequent releases -> monthly/bi-weekly
- Provenance support -> at fact level
  - same facts but different sources, basis for fusion
- Support for browsing data
- ID Management -> same entity in different Wikipedia editions
  - <u>http://global.dbpedia.org</u>
- Dataset dependency declaration
  - my service uses: dbpedia/Mapped-Geocoordinates/2017.05.10

Ш	ш	11	11	11	Ш	
Ŧ	Ħ.	Ŧ	Ŧ	Ŧ	Ŧ	

#	#	#	#	#####	#	#	###	###	#	#	#	###
#	#	#	#	#	#	#	#	#	#	#	#	
#	#	#	#	#	#	#	###	###	#	#	#	###
#	#	###	###	#	###	###	#	#	#	#		#
#	#	#	#	#	#	#	#	#	#	#	#	#
#####	#	#	#	#	#	#	###	###	#	###	#	###



**Digital Factory Platform** 

https://databus.dbpedia.org/

# Inspired by Maven Central

# **Databus - Digital Factory Platform**



https://databus.dbpedia.org

- Hosts a public metadata repository
- Not exclusive to DBpedia data
- Free to use, published metadata must be CC-0
- Published files stay on publisher's server
  - Full control over access (HTTP-Auth) and license

### **Automated monthly releases**



ABOUT

Classification of instances with the DBpedia Ontology. Contains triples of the form `<\$resource> rdf:type <\$dbpedia\_ontology\_class>` generated by the mappings extraction.

#### ARTIFACT VERSIONS

Version	Release Date	License
2019.10.01 (Documentation)	Oct 1st 2019	http://purl.oclc.org/NET/rdflicense/cc-by3.0
2019.09.01 (Documentation)	Sep 1st 2019	http://purl.oclc.org/NET/rdflicense/cc-by3.0
2019.08.30 (Documentation)	Aug 30th 2019	http://purl.oclc.org/NET/rdflicense/cc-by3.0
2019.08.01 (Documentation)	Aug 1st 2019	http://purl.oclc.org/NET/rdflicense/cc-by3.0
2019.07.01 (Documentation)	Jul 1st 2019	http://purl.oclc.org/NET/rdflicense/cc-by3.0
2019.06.01 (Documentation)	Jun 1st 2019	http://purl.oclc.org/NET/rdflicense/cc-by3.0
2019.05.01 (Documentation)	May 1st 2019	http://purl.oclc.org/NET/rdflicense/cc-by3.0
2019.04.20 (Documentation)	Apr 20th 2019	http://purl.oclc.org/NET/rdflicense/cc-by3.0
2019.04.01 (Documentation)	Apr 1st 2019	http://purl.oclc.org/NET/rdflicense/cc-by3.0
2019.03.01 (Documentation)	Mar 1st 2019	http://purl.oclc.org/NET/rdflicense/cc-by3.0
2019.01.01 (Documentation)	Jan 1st 2019	http://purl.oclc.org/NET/rdflicense/cc-by3.0



#### **Data Provenance and Versioning: DataID**

######

#	#	#	#	#####	#	#	###	##	#	#	##	##
#	#	#	#	#	#	#	#	#	#	#	#	
#	#	#	#	#	#	#	###	##	#	#	##	##
#	#	###	###	#	###	###	#	#	#	#		#
#	#	#	#	#	#	#	#	#	#	#	#	#
####	###	#	#	#	#	#	###	##	##	##	##	##

# Plugin version 1.3-SNAPSHOT - https://github.com/dbpedia/databus-maven-plugin

@prefix databus: <https://databus.dbpedia.org/> .
@prefix dataid-mt: <http://gataid.dbpedia.org/ns/mt#> .
@prefix dat: <http://gataid.dbpedia.org/ns/core#> .
@prefix dataid: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .
@prefix dataid-cv: <http://dataid.dbpedia.org/ns/core#> .
@prefix dataid-cv: <http://www.w3.org/2001/XMLSchema#> .
@prefix xds: <http://www.w3.org/2001/XMLSchema#> .
@prefix rdfs: <http://www.w3.org/ns/cat#> .
@prefix dataid-cv: <http://www.w3.org/ns/prow#> .
@prefix dataid-cv: <http://www.w3.org/ns/cat#> .
@prefix dataid-cv: <http://www.w3.org/ns/cat#> .
@prefix dataid-cv: <http://www.w3.org/ns/prow#> .
@prefix dataid-dv: <ht

<http://dbpedia-mappings.tib.eu/release/mappings/instance-types/2019.10.01/dataid.ttl#instance-types\_lang=el.ttl.bz2>

a	dataid:SingleFile ;
dataid:associatedAgent	<https: vehnem.github.io="" webid.ttl#this=""> ;</https:>
dataid:compression	"bzip2" ;
dataid:contentVariant	"el" ;
dataid:duplicates	"0"^^xsd:decimal ;
dataid:file	<https: 2019.10.01="" databus.dbpedia.org="" instance-types="" instance-types_lang="e1.ttl.bz2" mappings="" marvin="">;</https:>
dataid:formatExtension	"ttl";
dataid:isDistributionOf	<pre><http: 2019.10.01="" dataid.ttl#dataset="" dbpedia-mappings.tib.eu="" instance-types="" mappings="" release=""> ;</http:></pre>
dataid:nonEmptyLines	"102249"^^xsd:decimal;
dataid . preview	"# started 2019-10-08T13:21:28Z\n <http: el.dbpedia.org="" f3f»ťźťfťffft="" fÿf≵fif_ftf\$f0="" resource=""> <http: 196<="" td="" www.w3.org=""></http:></http:>

dataid:sha256sum "3ef05a493808099ffe9d90563aa712249eab325f424e6f485dc5893b117f00f2";

dataid:signature

"OxmICj7sQ8jIhzbGc9aohiBLoSa55poX4iUTid5bdjq/plk3519CdspZbyJD7+rvbTvs09JqO/+1sqW+/rCZe3eeqVz6+4XQbI8Uil0bTk5/3rShFQgjwmQo/aD9d/ry5C4xpyWWJHMb5qeSkwmyh6eeMc8WLoJSGqYVacik1wiXcUF9Pr770 cMYOoQHWMdwtIWoX0Gm07B+ltwosKQOCNpRuRK0+2iffN8Ashaz0BmynLgInilSQymZktf4oxDzLtyLeExCqPYsdjUbAuI4JebmrtGF6ciDOBuWSjH10LkAOfe/I2sG73yYCg6XgqhpBkmBJUcF95sxpIbDcj0blA==";

dataid:sorted false ; "16243387"^^xsd:decimal : dataid:uncompressedByteSize dataid-cv:lang "el" ; "http://dataid.dbpedia.org/ns/core#" ; dct:conformsTo dct:hasVersion "2019.10.01" ; dct:issued "2019-10-01T00:00:00Z"^^xsd:dateTime ; dct:license <http://purl.oclc.org/NET/rdflicense/cc-by3.0> ; dct:modified "2019-10-10T11:42:41Z"^^xsd:dateTime ;

# The Community



- Regular meetings
  - 14 community meetings since 2014 + several meetups
  - Usually co-located with SEMANTiCS or LDK conferences
  - o <u>https://wiki.dbpedia.org/join/community-meetings</u>
- Language chapters
  - German, Czech, Polish, Dutch, Italian, Catalan, Spanish, ...

#### Impact

- Over 26,000 DBpedia related papers
- 400 developers/researchers
- 600K file downloads per year
- 20 million hits daily (all APIs)
- 400 developers/researchers
- Inlinks by all major datasets
- 4000 website visitors weekly
- + and growing

=	Google Scholar	dbpedia Q
•	Articles	About 26,500 results
	Any time Since 2019 Since 2018 Since 2015 Custom range	Dbpedia: A nucleus for a web of open data           SAuer, C Bizer, G Kobilarov, J Lehmann, R Cyganiak The semantic web, 2007 - Springer           DBpedia is a community effort to extract structured information from Wikipedia and to make this information available on the Web. DBpedia allows you to ask sophisticated queries against datasets derived from Wikipedia and to link other datasets on the Web to Wikipedia           ☆ ワワ         Cited by 3930         Related articles         All 33 versions         ⊗
	Sort by relevance Sort by date	итиц DBpedia-A crystallization point for the Web of Data <u>C Bizer, J Lehmann</u> , G Kobilarov, <u>S Auer</u> - Web Semantics: science, 2009 - Elsevier
	<ul> <li>✓ include patents</li> <li>✓ include citations</li> </ul>	and to make this information accessible on the Web. The resulting <b>DBpedia</b> knowledge base currently describes over 2.6 million entities. For each of these entities, <b>DBpedia</b> defines a \$\frac{1}{2}\$ 90 Cited by 2273 Related articles All 18 versions Web of Science: 725
	Since Create alert	DBpedia spotlight: shedding light on the web of documents           PN Mendes. M Jakob, A García-Silva Proceedings of the 7th, 2011 - dl.acm.org           Interlinking text documents with Linked Open Data enables the Web of Data to be used as background knowledge within document-oriented applications such as search and faceted browsing. As a step towards interconnecting the Web to Documents with the Web of Data           ☆ 59         Cited by 1057         Related articles         All 15 versions         So
		DBpedia-a large-scale, multilingual knowledge base extracted from Wikipedia           J_Lehmann. R Isele, M Jakob, <u>A Jentzsch</u> Semantic, 2015 - content.iospress.com           The DBpedia community project extracts structured, multilingual knowledge from Wikipedia           and makes it freely available on the Web using Semantic Web and Linked Data           technologies. The project extracts knowledge from 111 different language editions of           ☆ 99 Cited by 1592 Related articles All 16 versions Web of Science: 361



## **Google Trends: DBpedia**





### **DBpedia Association members**









https://wiki.dbpedia.org - main website, general info

https://databus.dbpedia.org -> Databus interface, find data

https://forum.dbpedia.org -> DBpedia forum, Q&A

https://downloads.dbpedia.org/repo/lts/ -> data location, but better
use the databus interface

http://dbpedia.slack.com -> slack channel





# Part 1: Latest developments in DBpedia Part 2: Entityclassfier.eu NER



### Background

#### Entityclassifier.eu NER 🐉

Recognition, Classification, Linking of Salient Named Entities

- Project initiated in May 2012
- Developed during the LinkedTV EU FP7 project
- Some parts developed as part of the LOD2 project

#### A System for <u>Recognition</u>, <u>Classification</u> and <u>Linking</u>

of **Salient** Named Entities

### http://Entityclassifier.eu

#### Entityclassifier.eu NER 🐉

Recognition, Classification, Linking of Salient Named Entities

Extraction, Disambiguation and Classification of	f Entities and Named Entities
The Charles Bridge is a famous historic bridge that crosses the Vltava river in Prague, Czech Republic.         Hypernyms:	Request timeout (in seconds):       60       60         Language of the input text       Provenance of types         Participation       Provenance of types         English       THD         German       DBpedia         Dutch       Yago         THD knowledge base       Linked Hypernyms Dataset         Local Wikipedia mirror       Live Wikipedia         Entity Linking method       Entity Spotting method         SFI Based       POS Pattern based
#2: <u>ArchitecturalStructure</u> for entity disambiguated as <u>ornanes ornage</u> #2: <u>ArchitecturalStructure</u> for entity disambiguated as <u>Charles Bridge</u> #3: <u>Place</u> for entity disambiguated as <u>Charles Bridge</u> #4: <u>route of transportation</u> for entity disambiguated as <u>Charles Bridge</u> #5: place for entity disambiguated as Charles Bridge	Lucene (enhanced) Wikipedia Search All Voting Surface Form Similarity
#6: <u>infrastructure</u> for entity disambiguated as <u>Charles Bridge</u> #7: <u>Bridge</u> for entity disambiguated as <u>Charles Bridge</u> #8: <u>Bridge</u> for entity disambiguated as <u>Charles Bridge</u> #9: <u>RouteOffransportation</u> for entity disambiguated as <u>Charles Bridge</u> #0: <u>RouteOffransportation</u> for entity disambiguated as <u>Charles Bridge</u>	THD type filter     Types of entities to extract       DBpedia Ontology     ØNamed Entities       DBpedia instances     Common Entities       ØAll     All
	Force long entity linking

The <u>Charles Bridge</u> is a famous historic bridge that crosses the <u>VItava</u> river in <u>Prague</u>, <u>Czech Republic</u>.

# **Core Features**

Entityclassifier.eu NER 🐉

- Implemented in Java and GATE: <u>https://gate.ac.uk</u>
- Entity spotting -> JAPE grammars -> POS tags based
- Entity classification/typing -> JAPE grammars
  - Hypernym extraction of from the respective Wikipedia article
  - "Prague is the capital and largest city in the Czech Republic..."
  - Hypernum linked to DBpedia entity, <u>http://dbpedia.org/resource/Capital</u>
- Entity linking
  - Most Frequent Sense based -> Wikipedia Lucene index
  - Surface Form Similarity
  - Context based (DBpedia abstract as entity description)
  - ... several additional modifications/combinations

# **Core Features (cont.)**

Entityclassifier.eu NER 🐉

- Entity Salience
  - Extraction of the most important entities
  - Entities that play an important role in the document
- Machine learning
  - Training data via crowdsourcing (Figure Eight, ex Crowdflower)
    - most salient, less salient, not salient annotations
  - Experimented with several ML algorithms
  - Random tree forest best performing algorithm

# **Unique Features**

Entityclassifier.eu NER 🐉

- Realtime entity classification
  - Executed against live Wikipedia API
- Temporal type validity
  - Diego maradona: now "football manager", past "footballer"
  - Content is Wikipedia is regularly updated (for popular entities)
- Extraction of salient entities
  - Documents contain important and unimportant entities
  - Filter only important entities

# **Lessons learned**

Entityclassifier.eu NER 🐉

Recognition, Classification, Linking of Salient Named Entities

- Simple approaches perform relatively well
  - POS based recognition, most-frequent-sense based linking
- Provide an API
  - People would rather use the API then install the tool on local host
- Measure the impact
  - Collect data via API key request form
  - Jan 2015 to Aug 2017 = 59 API request, 42 unique users, 18 countries, 22 institutions, 59 English, 11 Dutch, 9 German
- Publish a paper with title the name of the tool (select keywords)

"Entityclassifier.eu: Real-Time Classification of Entities in Text with Wikipedia"





# Thank you

#### Milan Dojčinovski

