# JSONpedia

## Facilitating consumption of MediaWiki content

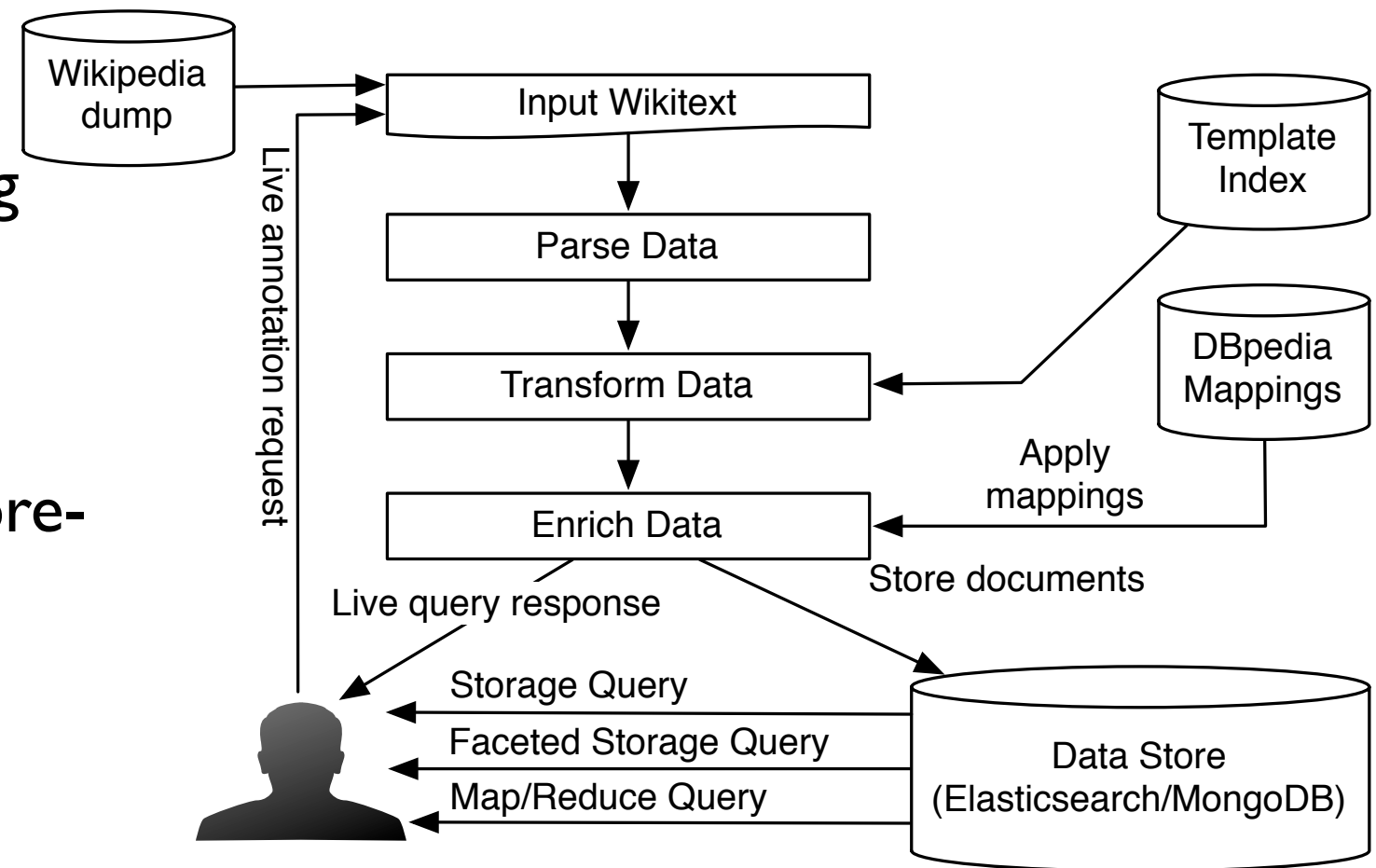*Michele Mostarda <mostarda@fbk.eu>, Twitter: @micmos*

v1.2

# Outline

- ‣ What is JSONpedia
- ‣ How does it work
- ‣ Main features
- ‣ Online demo
- ‣ Web App
- ‣ REST API
- ‣ jQuery plugin
- ‣ Code snippets
- ‣ Internals
- ‣ GSoC 2014
- ‣ History & previous work
- ‣ Forthcoming features
- ‣ Next release
- ‣ Online resources
- ‣ Support up
- ‣ Acknowledgements

# What is JSONpedia

JSONpedia is a Java library and a REST service meant to read MediaWiki pages as JSON.

# How does it work

- A user can perform a **live annotation** requests providing **Wikitext** or a reference to a Wikipedia page.
- A user can perform a **storage query** over the data storage pre-populated with the Wikipedia dump.
- A user can perform a **faceted storage query** over the data storage pre-populated with the Wikipedia dump.



- A user can perform a **faceted storage query** over the data storage pre-populated with the Wikipedia dump.
- A user can perform a **map/reduce storage query** over the data storage pre-populated with the Wikipedia dump.
- Any provided Wikitext is **parsed** (Parse Data), **templates are expanded** and new **metadata** is generated (Transform Data), **external data sources** are linked (Enrich Data), the final model is **converted in JSON** and stored into the Data Store.

# Main features

- ‣ WikiText event-based parser
- ‣ Configurable page processing pipeline
- ‣ Wikimedia template processing support
- ‣ DBpedia mapping integration
- ‣ RESTful interface
- ‣ MongoDB storage and map/reduce support
- ‣ Elasticsearch query support
- ‣ Elasticsearch faceting support
- ‣ Web frontend
- ‣ HTML data rendering
- ‣ CLI interface

# Online Demo

The official JSONpedia online demo is available at http://jsonpedia.org

# Web App

The JSONpedia web app allows to experiment with the REST service through a comfortable UX

# Live Panel

## Analyze any MediaWiki page live or directly copy/paste WikiText

# Query panel: MongoDB

## Query MediaWiki pages stored in MongoDB

### MongoDB Map/Reduce

GET /storage/mongo/mapred

Specify a data criteria selector, a map/reduce functions and optionally a resultset limit

**Criteria**  [<FIELD> OP <VALUE>]+    Criteria samples ▾    ?

[ Map/reduce function samples ▾ ]

**Map function**  function() { this... emit(key, val); }    ?

**Reduce function**  function(key, vals) { return ... }    ?

**Limit**  1000    ?

Query ☑

Run    Cancel

### MongoDB Query

GET /storage/mongo/select

Specify a data selector and optionally a filter and a resultset limit

**Selector**  [<FIELD> OP <VALUE>]+ -> [<FIELD>]+    Selector samples ▾    ?

**Filter**  [<FIELD> : <VALUE re>]+    Filter samples ▾    ?

**Limit**  1000    ?

Query ☑

Query    Cancel

# Query panel: Elasticsearch

Query the latest Wikipedia dump with **Elasticsearch**

## Elasticsearch Query

GET /storage/elastic/select

Specify a data selector and optionally a filter and a resultset limit

| | | |
|---|---|---|
| **Selector** | [[<FIELD>:]?<CRITERIA>]+ | Selector samples ▾ |
| **Filter** | [<FIELD> : <VALUE re>]+ | Filter samples ▾ |
| **Limit** | 1000 | |

Query ✎

Query    Cancel

# Query panel: Elasticsearch

Explore the latest Wikipedia dump with
**Elasticsearch** *FacetView*

# REST API

**GET** `/annotate/resource/{json|html}/{res-id|res-url}`

Process a live WikiMedia resource

**POST** `/annotate/resource`

(wikitext, format, processors, filter)
Process arbitrary WikiText markup

**GET** `/storage/mongo/select`
`?q=<query>&filter=<filter>&limit=<limit>`

Query the Wikipedia dump with MongoDB

**GET** `/storage/mongo/mapred`
`?map=<map-func>&red=<red-func>&criteria=<criteria-exp>&limit=<limit>`

Query the Wikipedia dump with MongoDB Map / Reduce

**GET** `/storage/elastic/select`
`?q=<query>&filter=<filter>&limit=<limit>`

Query the Wikipedia dump with Elasticsearch

# jQuery Plugin

JSONpedia comes with a jQuery 1.8 plugin providing facilitated access to the REST service.

http://jsonpedia.org/frontend/js/jsonpedia.js

# Code Snippets

*Example:*
*retrieve content of page London from English Wikipedia, extract the DOM structure, filter nodes of type "section", get first of them and render as HTML.*

```java
import com.machinelinking.main.JSONpedia;
import org.codehaus.jackson.JsonNode;

JSONpedia jsonpedia = JSONpedia.instance();
JsonNode root = jsonpedia.process("en:London").flags("Structure").json();

JsonNode[] sections = jsonpedia.applyFilter("@type:section", root);
String firstSectionHTML = jsonpedia.render("en:London", sections[0]);
```

# Internals

# Processing Pipeline



This picture shows the processing pipeline implemented in JSONpedia

# Types of Processor

A Processor receives a stream of events generated by parser and perform data enrichment and transformation.

‣ Structure

‣ Extractors

‣ Linkers

‣ Splitters

‣ Validator

# Structure

The *Structure* Processor receives a stream of WikiText parsing events and builds a 1-1 JSON representation of the document DOM.

# Extract

Extractors are specific Processors that collect a certain type of data from the event stream.

*For example the SectionsExtractor collects a list of all sections declared in the document stream*

# Split

A *Splitter* is a Processor cutting sub-trees of the JSON document built by the Structure processor.

*An example of Splitter is the TableSplitter which collects the JSON nodes representing all tables found in document.*

# Link

A *Linker* is a Processor which links the detected document entities to other information acquired from external sources.

*An example of Linker is the FreebaseLinker which connects an entity to the same representation in Freebase if any.*

# Validate

A *Validator* is a Processor performing the check of data structures parsed from a document.

# WikiText event based parser messages

```java
// Document bounding.
void beginDocument(URL document);
void endDocument();

// Error handling.
void parseWarning(String msg,
ParserLocation location);
void parseError(Exception e,
ParserLocation location);

// Tag handling.
void beginTag(String node, Attribute[]
attributes);
void endTag(String node);
void inlineTag(String node,
Attribute[] attributes);
void commentTag(String comment);

// Sections
void section(String title, int level);

// References
void beginReference(String label);
void endReference(String label);

// Links
void beginLink(String url);
void endLink(String url);

// lists
void beginList();
void listItem();
void endList();

// Templates
void beginTemplate(String name);
void endTemplate(String name);

// Tables
void beginTable();
void headCell(int row, int col);
void bodyCell(int row, int col);
void endTable();

// Generic parameter
void parameter(String param);
// Plain text
void text(String content);
```

# JSONpedia @Google Summer of Code 2014

**Project**:

JSONpedia Extractor

**Organization**:

DBpedia & DBpedia Spotlight

**Student**:

Roberto Bampi

**Mentor**:

Michele Mostarda

**Description**:

Create a general infrastructure to create DBpedia extractors based on JSONpedia.

**Public Repo**:

https://github.com/dbpedia/jsonpedia-extractor/

The JSONpedia extractor for DBpedia relies on a Wikipedia dump processed with JSONpedia and stored in Elasticsearch, and allows to build scriptable data scrapers based on faceted queries.

# Extraction Samples

**Discography**
Extract artist, album, year and reference for all discographies defined in Wikipedia

**Painter works**
Extract painter, work, year and link for any paining defined in Wikipedia

**Public Gardens**
Extract city, garden, description for any public garden defined in Wikipedia

# Forthcoming Features

‣ JSONpedia dumps will be available for download.

‣ RDF output.

‣ Online data model Exporter Tool (CSV).

*Follow the updates here: https://bitbucket.org/hardest/jsonpedia/issues*

# History & Previous Work

# History

‣ Initially conceived as a tool to generate machine learning datasets.

‣ The REST service,inspired by Sweeble Crystalball,produces JSON and a browsable HTML data.

‣ Written over a context-dependent event based parser to be more performant than a regex matcher (like the WikiParser) or a DOM based parser (like Sweeble).

# Differences with DBpedia

‣ JSONpedia produces JSON, DBpedia RDF.

‣ JSONpedia includes all the structural elements of a page: links, references, lists, sections, template, tables, XML markup.

‣ JSONpedia produces low-refined data which requires further processing to be consumed, DBpedia produces ready to use high quality data.

‣ JSONpedia is a not competitor of DBpedia but rather a complement.

# Differences with Sweeble

‣ Lightweight Event based parser vs DOM parser.

‣ More tolerant to frequent syntax errors present within WikiText pages.

‣ Serializes to JSON output which is easier to consume!

# Next Release

*End of March 2015 v1.2*

# Online resources

live demo:

http://jsonpedia.org/

source code:

https://bitbucket.org/hardest/jsonpedia

# Acknowledgements

**Marco Fossati** - FBK WeD, PhD student, DBpedia community member.

*@hjfocs*

**Roberto Bampi** - SpazioDati, Backend Developer, JSONpedia contributor and student in GSoC 2014.

*@BampiRoberto*

# Thanks for reading!

*Michele Mostarda <mostarda@fbk.eu>, Twitter: @micmos*