

Publishing Linked Data as Datasets

CS-E4410 Semantic Web, 26.1.2022

Eero Hyvönen

*Aalto University, Semantic Computing Research Group (SeCo) <http://seco.cs.aalto.fi>
University of Helsinki, HELDIG <http://heldig.fi>*

eero.hyvonen@aalto.fi

Learning Objective

Understand how Linked Data is published on the Web as datasets and SPARQL endpoints to be used for practical applications

Outline

Web of Data: Basic principles

Distributed Approach: Embedding data in HTML pages

Centralized Approach: Standalone data services on servers

Web of Data: Basic Concepts

Basic Concepts Clarified

Linked Data

- Means practical and simple Web of Data using RDF graphs
 - Semantic Web technologies include also more advanced models
- Based on W3C's Semantic Web standards
 - *"Rebranding Semantic Web", focus on simple practical semantics*
- Linked Data is often open but can be closed, too

Open Data

- Openly available data (on the Web), under an open license
 - Cf. Creative Commons Licenses: <https://creativecommons.org/>
 - *"Open data and content can be freely used, modified, and shared by anyone for any purpose" – <http://opendefinition.org>*
- Open data is not necessarily free of charge

What is Web of Data?

The Web we see is a network that links pages: Web of Pages

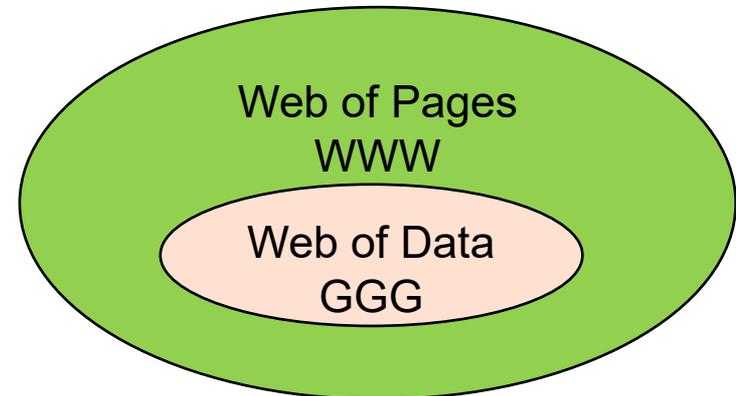
- Can be searched/browsed with a traditional web browser
- Links connect pages

A data network has emerged inside the web: Web of Data

- Can be searched/browsed with a semantic application
- Links in Web of Data connect concepts and data (e.g., eagle → bird)

Semantic Web consists of *both* networks

- Web of Pages (for humans)
 - *WWW World Wide Web*
- Web of Data (for machines to use)
 - *GGG Giant Global Graph*
 - *Examples of knowledge graphs in applications:*
 - Google Knowledge Graph, Microsoft Satori, Facebook Open Graph, ...
 - Lots of domain specific applications

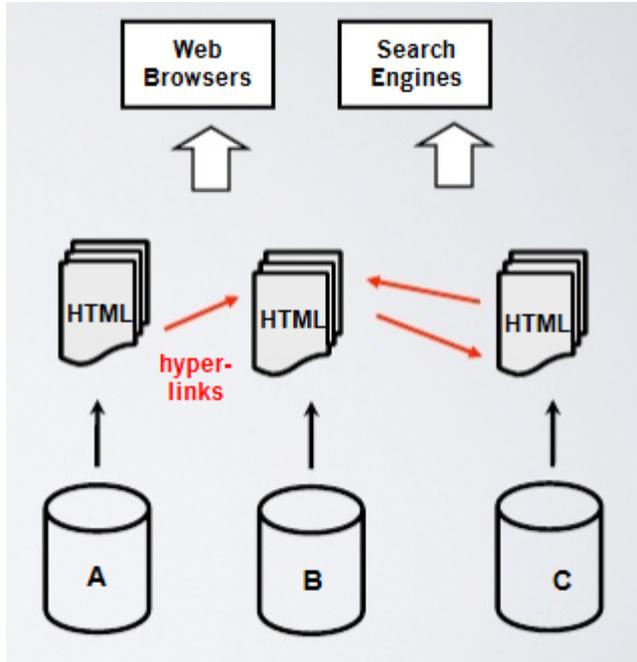


Idea of Linked Data

- Datasets are created and published using RDF
 - Embedded in web pages or as data services / data bases
 - Lightweight Semantic Web technologies are used
 - “A Little Semantics Goes a Long Way“ – Jim Hendler
- Datasets are linked together to enrich their content
 - *Cross-referencing data in other datasets*
 - E.g., place ”Finland” in GeoNames.org to president ”Niinistö” in DBpedia.org
 - *Identifying same concepts in different datasets (data alignment)*
 - E.g., ”Helsinki” in GeoNames.org vs. DBpedia.org
- Enriched data is (re-)published as data services for applications
- Coordinated by Linked Open Data (LOD) communities
 - *Application domain specific and global ones*

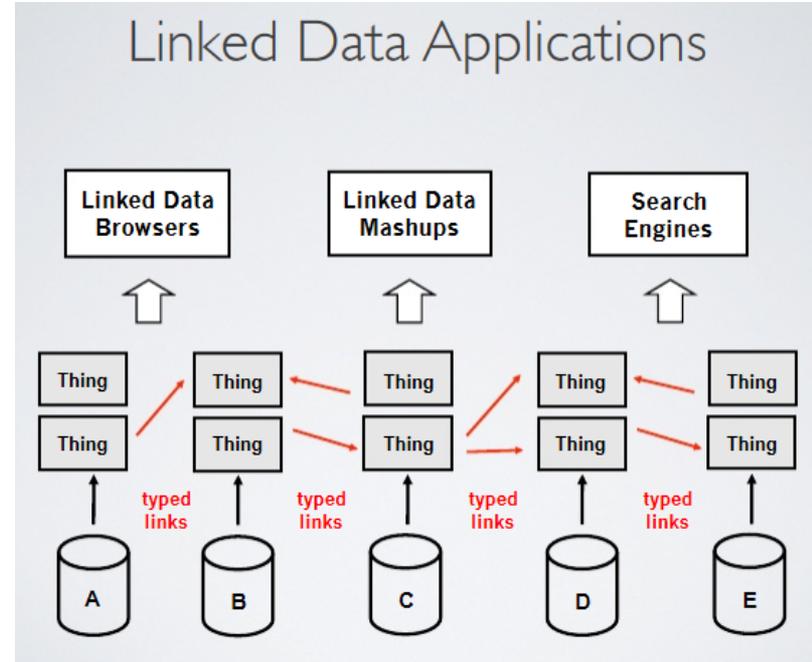
WWW and GGG coexistence

WWW



(Anja Jentzsch, 2012)

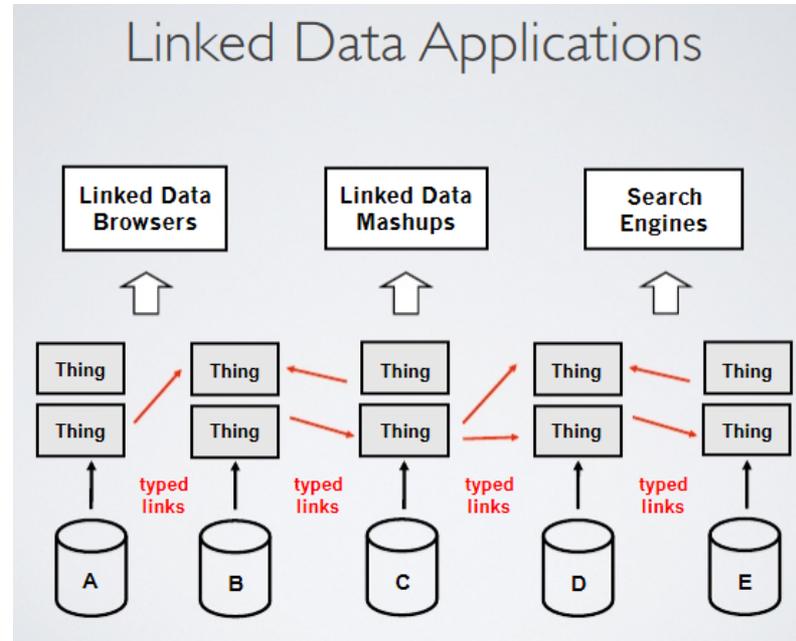
GGG



(Anja Jentzsch, 2012)

Dimensions of application development

- Services for
 - *Humans*
 - *Machines*
- Data linking
- Data aggregation
- Data harvesting
- Data production



(Anja Jentzsch, 2012)

Publishing Linked Data as Datasets for the Web of Data

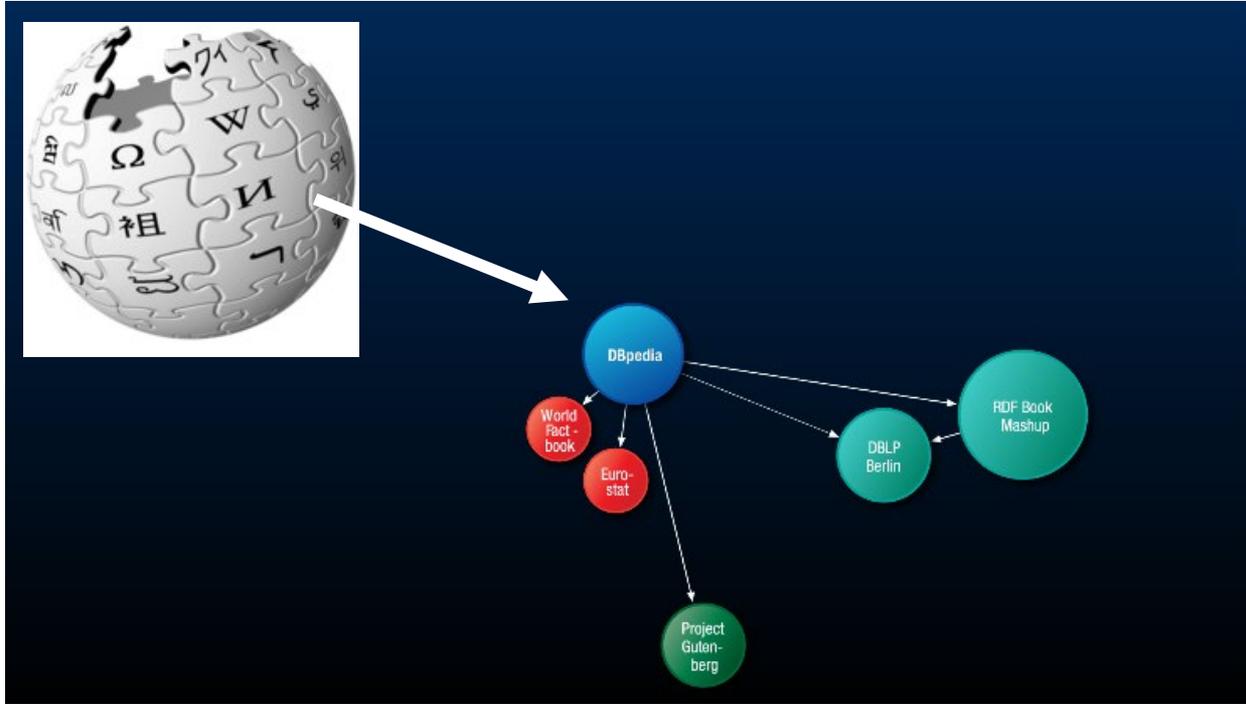
Dataset Publishing: two ways

Publish datasets for downloading or via APIs

- **CSV files**
- **RDF files**
- **JSON files**

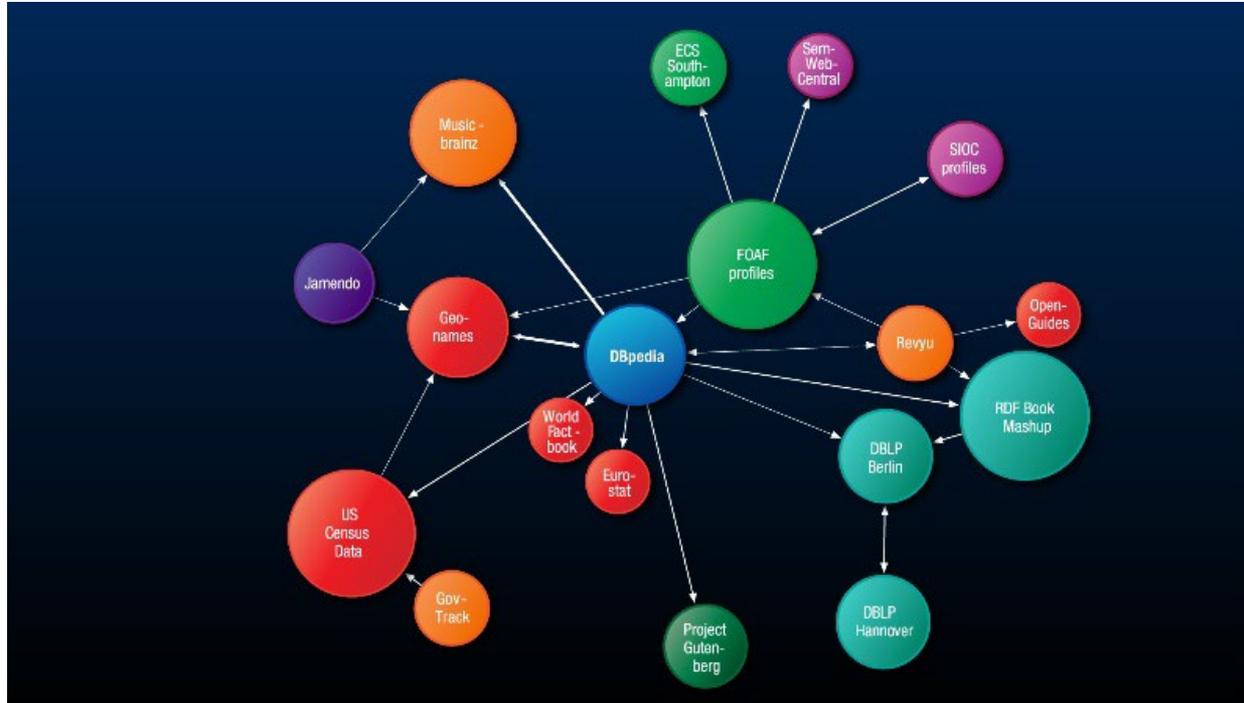
Publish Linked Data in SPARQL endpoints ready for querying

Publishing Linked Open Datasets: Story Starting 2005



(Tim Berners-Lee)

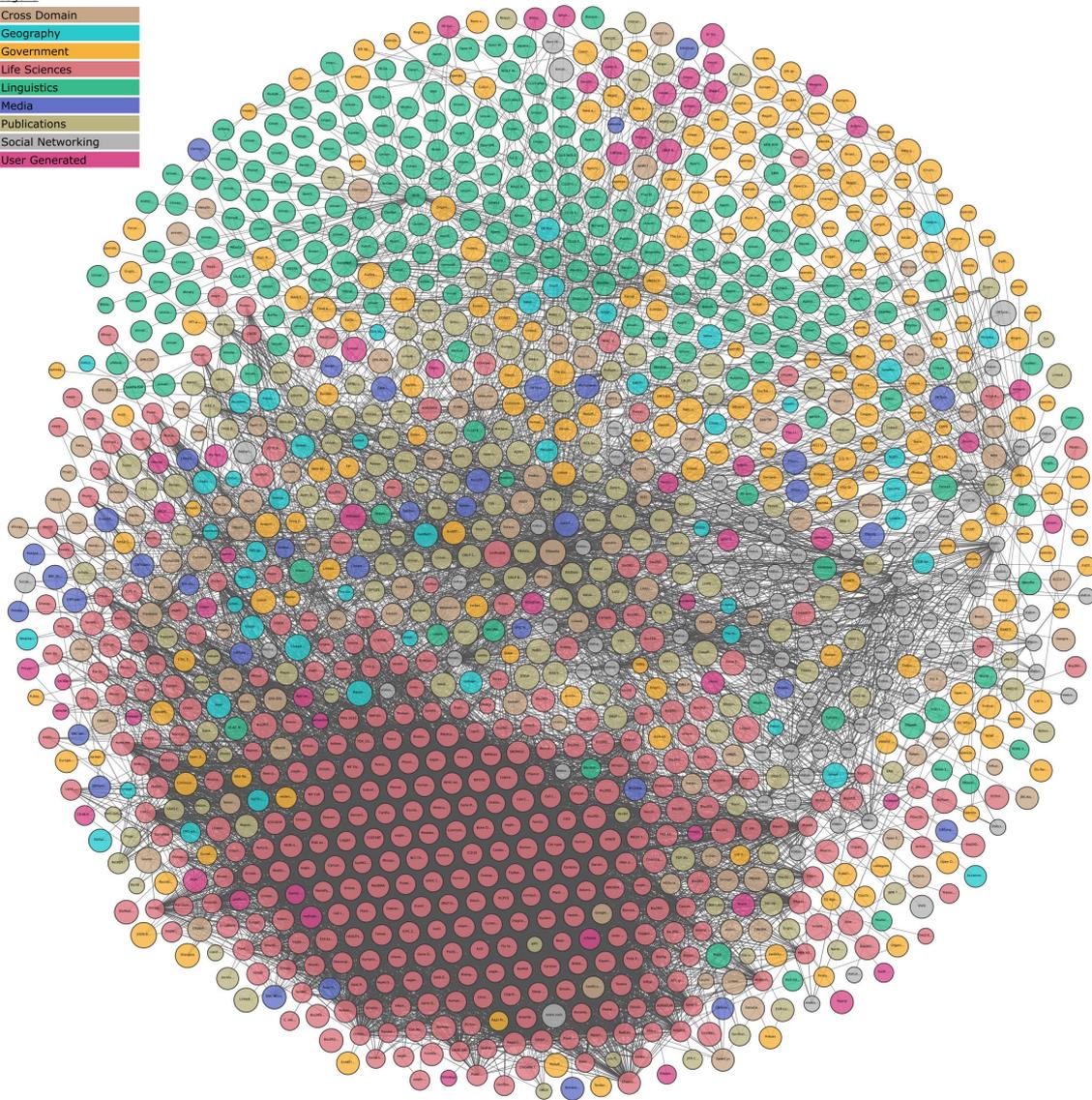
New Datasets Linked with Others



(Tim Berners-Lee)

2021

<https://lod-cloud.net/>



There are SPARQL endpoints around

<https://www.w3.org/wiki/SparqlEndpoints>



Page Discussion

Read View source View history Search

SparqlEndpoints

see also: SPARQL

In addition to the list below, Mondeca provides a [SPARQL endpoint uptime service](#) which monitors the availability of all SPARQL endpoints that are cataloged in [CKAN](#). A [similar service](#) is provided by Vienna University.

Currently Alive SPARQL Endpoints

(alphabetical. let's avoid [PoorMansHypertext](#) and in-your-face URIs, please)

Project	status	SPARQL endpoint	Webform	comment
Wikidata	(2017-02-23) alive	endpoint	GUI	See also SPARQL federation input
BBC Programmes and Music	(2010-06-29) alive	endpoint	Ajax based Visual Query Builder	Powered by OpenLink Virtuoso ; also supports Faceted Browsing and Exploration
Bio2RDF	(2010-01-07) alive	List of 40 SPARQL endpoints	n/a	uses OpenLink Virtuoso
BioGateway	(2010-01-07) timeout	endpoint	webform	BioGateway provides many parameterizable SPARQL queries, both biological as ontological, on RDF graphs that were optimized for querying. The graphs have relational closures. Empowered by OpenLink Virtuoso .
BBC Backstage (HP Labs)	(2010-01-07) server not responding	endpoint	webform	uses joseki 3
BBC John Peel sessions from DBTune (Centre for Digital Music, Queen Mary, University of London)	(2010-01-07) alive	endpoint	n/a	dbtune aims to gather and interlink music-related information.
BBC playcount data from DBTune (Centre for Digital Music, Queen Mary, University of London)	(2010-01-07) alive	endpoint	n/a	dbtune aims to gather and interlink music-related information.
DailyMed	(2010-01-07) alive	endpoint		a D2R endpoint
data.gov	(2010-05-22) alive	endpoint	webform	uses OpenLink Virtuoso
data.gov.uk	(2010-02-04) alive	endpoint	webform	The data.gov.uk endpoint
DBLP Bibliography Database published through D2R Server (Freie Universität Berlin)	(2010-01-07) alive	endpoint	webform (Maybe only Firefox)	The DBLP database provides bibliographic information on major computer science journals and conference proceedings.
DBpedia (University of Mannheim, Universität Leipzig, OpenLink Software)	(2010-01-07) alive	endpoint	SNORQL.webform (Firefox/Safari/Opera); Ajax based Visual Query Builder	dbpedia.org is a community effort to extract structured information from periodic Wikipedia dumps and to make this information available on the Web. It is served to the public via a live instance of OpenLink Virtuoso , and also offers Faceted Browsing and Exploration
DBpedia-live (Universität Leipzig, University of Mannheim, OpenLink Software)	(2010-01-07) alive	endpoint	webform	Based on, now parallel to, and soon to replace the existing dbpedia.org data sets, DBpedia-Live is constantly updated, based on Wikipedia change-feeds. It is served to the public via a live instance of OpenLink Virtuoso , and also offers Faceted Browsing and Exploration
German DBpedia (AG Corporate Semantic Web, Freie Universität Berlin)	(2008-10-15) alive	endpoint	site	de.dbpedia.org is the German language chapter of DBpedia
DBpedia Live German (AG Corporate Semantic Web, Freie Universität Berlin)	(2008-10-15) alive	endpoint	site	de.dbpedia.org is the German language chapter of DBpedia
Spanish DBpedia (Universidad Autónoma de Madrid, Universidad Politécnica de Madrid, OpenLink Software)	(2011-04-04) alive	endpoint	site	es.dbpedia.org is the Spanish chapter of DBpedia
Diseasome	(2010-01-07) alive	endpoint		a D2R endpoint
DoapSpace	(2010-01-07) away	endpoint	webform	This is a highly experimental TurboGears with rdflib triplestore (mysql) SPARQL endpoint.
DrugBank	(2010-01-07) alive	endpoint		a D2R endpoint
EEA (European Environment Agency) Semantic	(2010-01-07) alive	endpoint		

- Main Page
- Browse categories
- Recent changes
- Tools
- What links here
- Related changes
- Special pages
- Printable version
- Permanent link
- Page information

An example of an RDF production pipeline: Case WarSampo

8

M. Koho et al. / WarSampo Knowledge Graph: Finland in WW2 as Linked Open Data

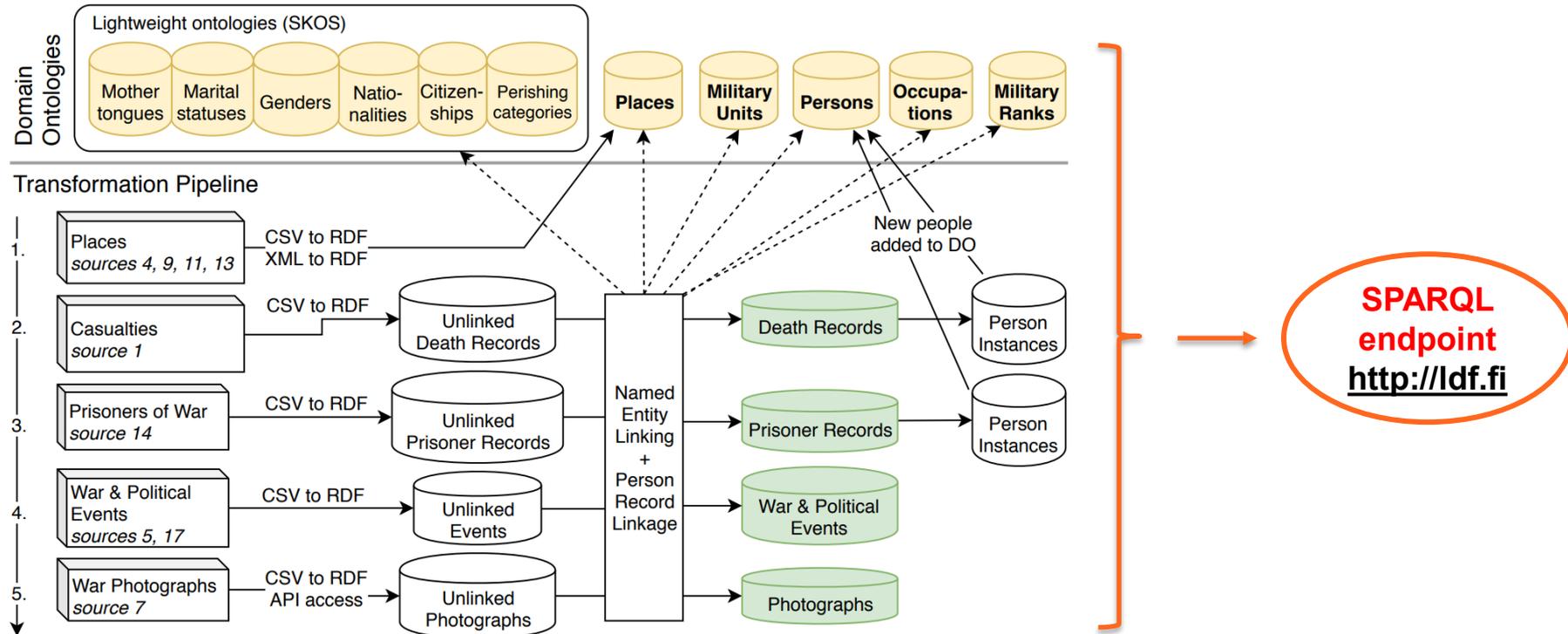
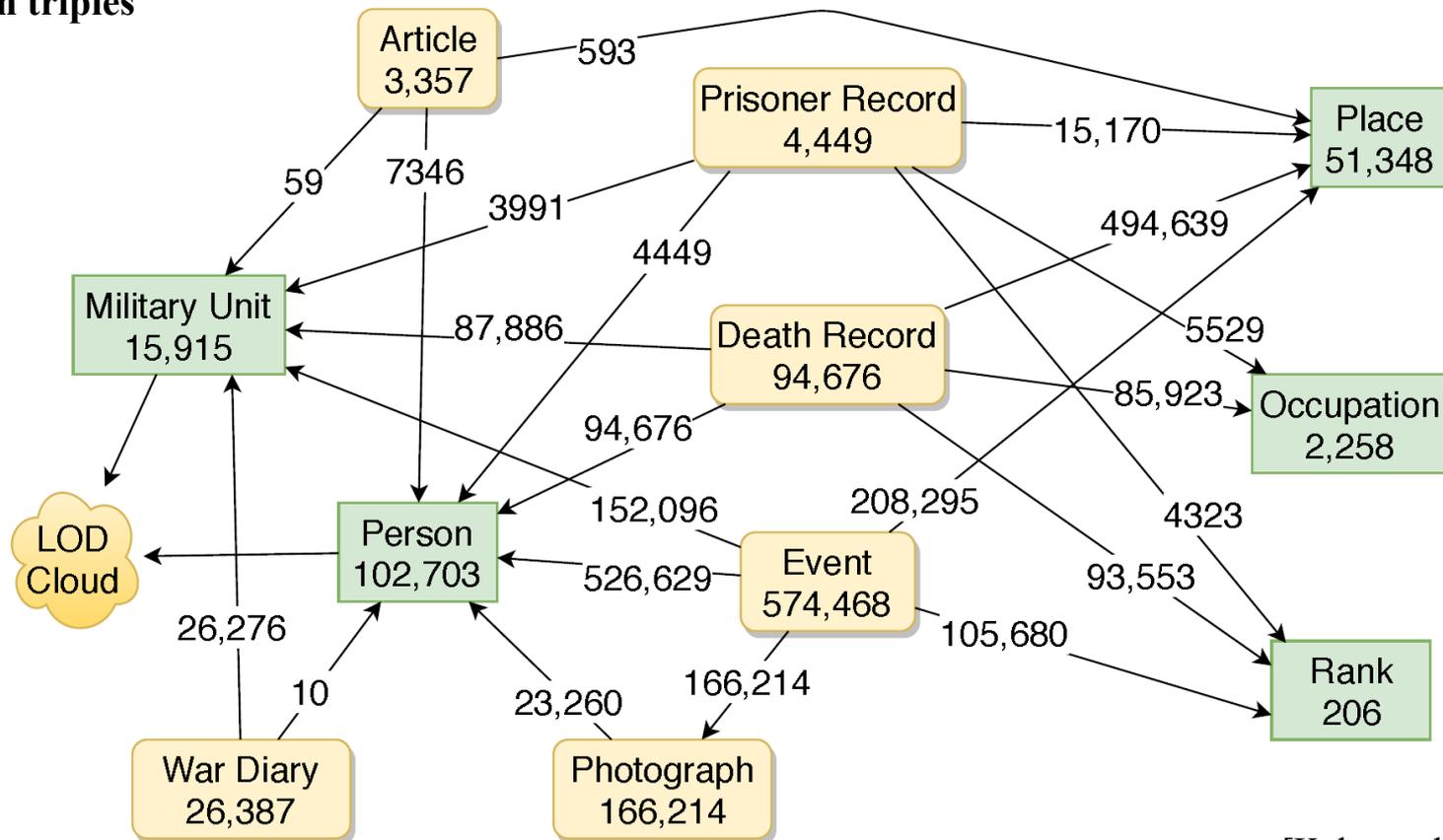


Figure 3. The 5-step WarSampo data transformation process. Dashed arrows represent entity linking, while solid arrows convey data flow.

WarSampo: A Linked Open Data cloud of its own

14 million triples

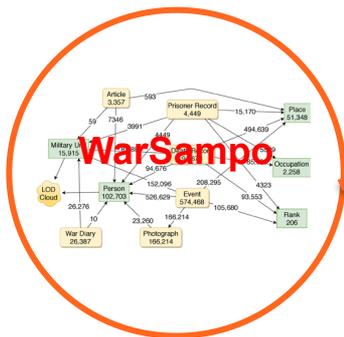
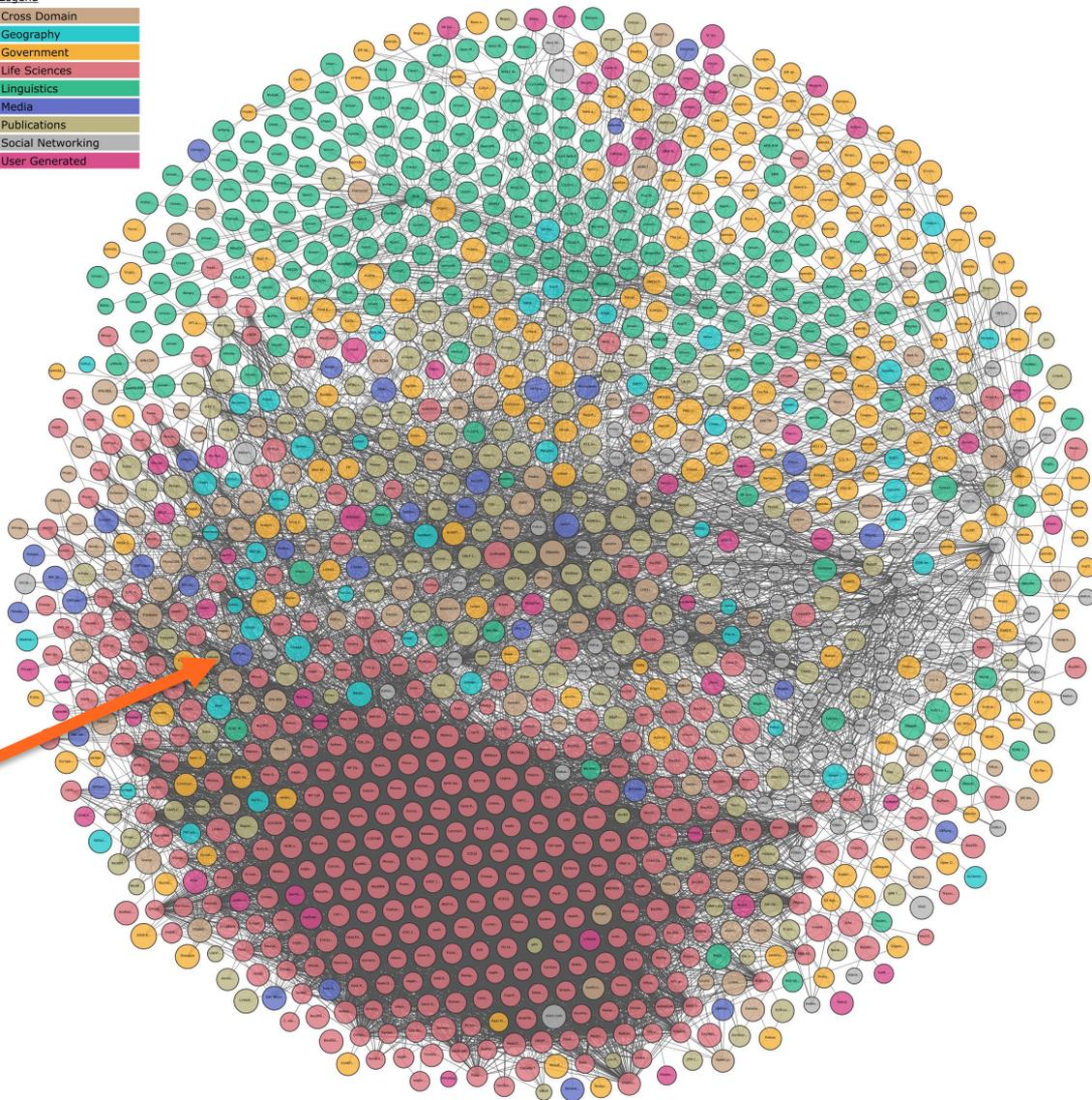


[Koho et al., 2021]

2021

<https://lod-cloud.net/>

- Legend
- Cross Domain
 - Geography
 - Government
 - Life Sciences
 - Linguistics
 - Media
 - Publications
 - Social Networking
 - User Generated



WarSampo dataset page & SPARQL endpoint: <http://www.ldf.fi/dataset/warsa>



[Home](#)
[Project](#)
[Datasets](#)
[Search Data](#)
[Schemas](#)
[Services](#)
[Policies](#)
[Documentation](#)
[Validation](#)
[Linked Data Science](#)
[Applications](#)
[Your Data?](#)
[Linked Data School](#)

WarSampo

Sotasampo

Linked Data Finland

★★★★★

WarSampo Knowledge Graph includes harmonized data of different kinds concerning the Second World War in Finland, separated in different subgraphs representing events, actors, places, photographs, and other aspects and documentation of the war. The data covers the Winter War 1939-1940 against the Soviet attack, the Continuation War 1941-1944 where the occupied areas of the Winter War were temporarily regained, and the Lapland War 1944-1945, where the Finns pushed the German troops away from Lapland.

To test and demonstrate its usefulness, this Knowledge Graph is in use in the semantic portal [WarSampo](#), explained in more detail in the [project page](#).

The Knowledge Graph is published on [Zenodo](#) with a version history

Example SPARQL queries for the data:

- [Events, photographs and articles that are situated in Vyborg](#)
- [Casualties of the 1st Division and its subunits in the time interval 13.2.-13.3.1940 by place and date](#)

Data Download

The data can be downloaded at <https://zenodo.org/record/3431122/files/warsampo.zip>.

License

CC BY 4.0

Licensors: [Kansallisarkisto](#), [Semanttisen laskennan tutkimusryhmä \(SeCo\)](#)

See possible graph-specific licenses below.

Detailed Dataset Contents

Karelian map names 1922-44 (URI: http://df.fi/warsa/places/karelian_places)



Information about Linked Data Available and Use Cases

Use cases

W3C's Working Group Note: Data on the Web Best Practices Use Cases & Requirements

1. Introduction
2. Use Cases
 - 2.1 ASO: Airborne Snow Observatory
 - 2.2 BBC
 - 2.3 Bio2RDF
 - 2.4 BuildingEye: SME use of public data
 - 2.5 Dados.gov.br
 - 2.6 Digital archiving of Linked Data
 - 2.7 Dutch Base Registers
 - 2.8 GS1 Digital
 - 2.9 ISO GEO Story
 - 2.10 The Land Portal
 - 2.11 LA Times' Reporting of Ron Galperin's Infographic
 - 2.12 LusTRE: Linked Thesaurus fRamework for Environment
 - 2.13 Machine-readability and Interoperability of Licenses
 - 2.14 Mass Spectrometry Imaging (MSI)
 - 2.15 OKFN Transport WG
 - 2.16 Open City Data Pipeline
 - 2.17 Open Experimental Field Studies
 - 2.18 Resource Discovery for Extreme Scale Collaboration (RDESC)
 - 2.19 Recife Open Data Portal
 - 2.20 Retrato da Violência (Violence Map)
 - 2.21 Share-PSI 2.0: Uses of Open Data Within Government for Innovation and Efficiency
 - 2.22 Tabulae - how to get value out of data
 - 2.23 UK Open Research Data Forum
 - 2.24 Uruguay Open Data Catalog
 - 2.25 Web Observatory
 - 2.26 Wind Characterization Scientific Study
3. General Challenges
 - 3.1 A Word on Open and Closed Data
 - 3.2 Requirements by Challenge

In the USA:

DATA.CATALOG

Search Data.Gov

DATA TOPICS - RESOURCES STRATEGY DEVELOPERS CONTACT

Organizations / Datasets

Search datasets...

Order by: Popular

Datasets ordered by Popular

Formats: RDF x

Filter by location

Enter location...

12,018 datasets found

Population Demographics, 1995-2012 [12,208 recent views](#)

State of Illinois - The number of induced pregnancy terminations reported in Illinois by county (if in excess of 50), by age and marital status. Note: Marital status and age are only...

CSV RDF JSON XML

State

Governor's Children's Cabinet County Crime Rates And Population [16 recent views](#)

State of Illinois - This dataset was compiled by the Illinois Criminal Justice Information Authority (ICJIA) at the request of the Governor's Children's Cabinet. This data contains the...

CSV RDF JSON XML

State

Most Popular Baby Boy Names, 1980-2013 [13 recent views](#)

State of Illinois - Note: 2010-2013 should be considered provisional data.

CSV RDF JSON XML

State

Youth Suicide Deaths in Washington State by Gender Age 0-17 Years, from 2008-2012 [11 recent views](#)

State of Washington - Youth Suicide Deaths in Washington State by Gender Age 0-17 Years, from 2008-2012

CSV RDF JSON XML

State

IDPH Population Projections For Chicago By Age And Sex 2010 To 2025

State of Illinois - Introduction This report presents projections of population from 2015 to 2025 by age and sex for Illinois, Chicago and Illinois counties produced for the Certificate...

CSV RDF JSON XML

State

Police Homeless-related Incidents

City of Santa Rosa - Homeless-related Police Incidents in Santa Rosa are identified and counted by any incident relating from a homeless-related call for service.

CSV RDF JSON XML

City

Leading Causes of Death, 1990-2010

State of Illinois - Leading causes of death for Illinois residents was compiled by the U. Center for Health

State

Map tiles & Data by OpenStreetMap, under CC BY-SA

Topics

Local Government (10279)

Older Adults Health... (77)

Climate (6)

Disasters (3)

Topic Categories

Human Health (6)

Dataset Type

non-geospatial (12016)

geospatial (2)

Tags

https://catalog.data.gov/dataset?res_format=RDF&res_format_limit=0

How to create a SPARQL endpoint and applications on top of it by yourself?

Apache Jena Fuseki

<https://jena.apache.org/documentation/fuseki2/>

[Apache Jena](#) [Home](#) [Download](#) [Learn](#) [Javadoc](#) [Ask](#) [Get involved](#) [Edit this page](#)

DOCUMENTATION / FUSEKI2

Apache Jena Fuseki

Apache Jena Fuseki is a SPARQL server. It can run as an operating system service, as a Java web application (WAR file), and as a standalone server.

Fuseki comes in in two forms, a single system "webapp", combined with a UI for admin and query, and as "main", a server suitable to run as part of a larger deployment, including [with Docker](#) or running embedded. Both forms use the same core protocol engine and [same configuration file format](#).

Fuseki provides the SPARQL 1.1 [protocols for query and update](#) as well as the [SPARQL Graph Store protocol](#).

Fuseki is tightly integrated with [TDB](#) to provide a robust, transactional persistent storage layer, and incorporates [Jena text query](#).

Contents

- [Download with UI](#)
- [Getting Started](#)
- [Running Fuseki with UI](#)
 - [As a standalone server with UI](#)
 - [As a service](#)
 - [As a web application](#)
 - [Security with Apache Shiro](#)
- [Running Fuseki Server](#)
 - [Setup](#)
 - [As a Docker container](#)
 - [As an embedded SPARQL server](#)
 - [Security and data access control](#)
 - [Logging](#)
- [Fuseki Configuration](#)
- [Server Statistics and Metrics](#)
- [How to Contribute](#)
- [Client access](#)
 - [Use from Java](#)
 - [SPARQL Over HTTP](#) - scripts to help with data management.
- [Links to Standards](#)

The Jena users mailing is the place to get help with Fuseki.

[Email support lists](#)

Download Fuseki with UI

Releases of Apache Jena Fuseki can be downloaded from one of the mirror sites:

[Jena Downloads](#)

and previous releases are available from [the archive](#). We strongly recommend that users use the latest official Apache releases of Jena Fuseki in preference to any older versions.

Fuseki download files

INSTRUCTIONS OF USING FUSEKI:

1. Download Fuseki and install it on your machine
2. Start Fuseki
3. You can use it at localhost with your browser
4. Upload an RDF file in the SPARQL endpoint
5. You are ready to query and develop applications!
6. The application can later be published on a web server



More Information – Questions?

Semantic Web & Linked Data Standards

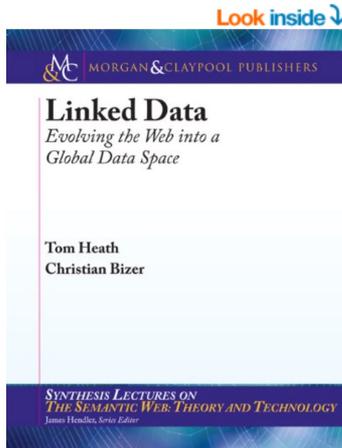
<http://www.w3.org/standards/semanticweb/>

T. Heath, C. Bizer: **Evolving the Web into a Global Data Space**

Free online version: <http://linkeddatabook.com/editions/1.0/>

WarSampo project homepage:

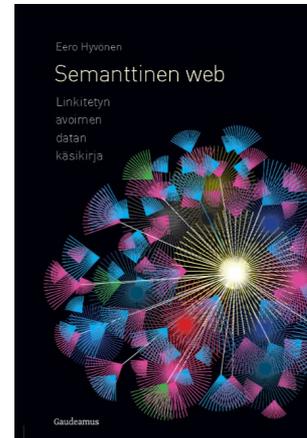
<https://seco.cs.aalto.fi/projects/sotasampo/en/>



In English

2011

<https://www.amazon.com/Linked-Data-Evolving-Global-Space-ebook/dp/B009KC1YM2>



In Finnish

2018

<https://www.gaudeamus.fi/semanttinen-web/>