# Publishing data on the Web

Bernadette Farias Lóscio and Caroline Burle

# Topics to be discussed

- Data on the Web

- Linked Data

- Open Data

- Big Data, Open Data and Data on the Web

- Data on the Web lifecycle

- DWBP: Challenges and Benefits

- Questions and comments

# Data on the Web
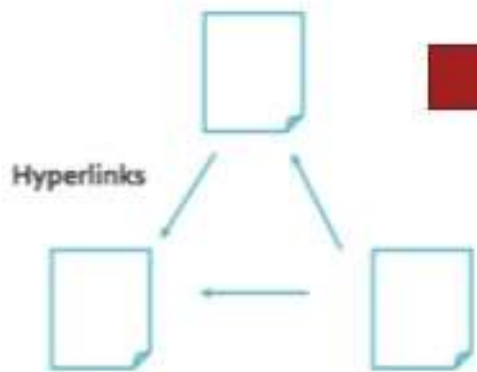


Consumes data

Publishing data on the web

Publishes data

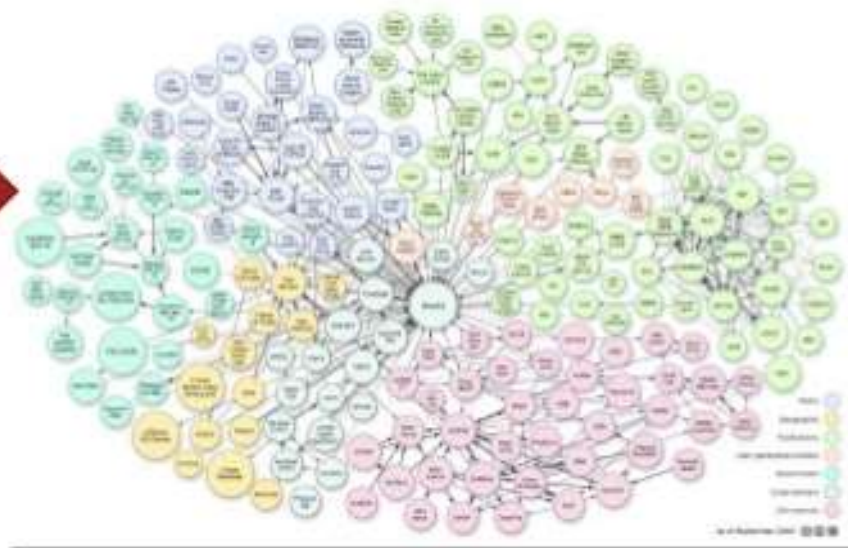# Web of Linked Documents → Web of Linked Data

## The Web is evolving from a "Web of linked documents" into a "Web of linked data"... (1/2)

Web of documents...

Web of linked data...

Hyperlinks

# What is Open Data?

Characteristics:

**Availability and access**

**Reuse and Redistribution**

**Universal Participation**



Source: http://ceweb.br/publicacao/open-data-guideline/

ceweb.br nic.br cgi.br

# De Olho nas Metas



Source: http://deolhonasmetas.org.br

# De Olho nas Metas



Source: http://deolhonasmetas.org.br/open-data

# 5-Star Open Data



Tim Berners-Lee, the inventor of the Web and Linked Data initiator, suggested a 5-star deployment scheme for Open Data. Here, we give examples for each step of the stars and explain costs and benefits that come along with it.

Source: http://5stardata.info/

**Big Data**

Publishing data

Publishing data on the Web

Open Data Principles

Data on the Web Best Practices

*The more stars have the data, the easier it is to find and reuse the data.*

ceweb.br  nic.br  cgi.br

How to make data available?

Which data to publish?

How to make data interoperable?

Publishing data on the Web

Which are the data sources?

How to identify data resources?

Which data formats to use?

How to gather feedback?

*Publishing data on the Web is more than just publishing!*

# Data on the Web Lifecycle



Planning

*Choice of which data to publish*

Feedback

Refining

Creation

Consuming

Publication

Web

Access

# Data on the Web Lifecycle



Diagram showing the Data on the Web Lifecycle with boxes: Planning, Feedback, Creation (circled in red), Publication, Web (cloud), Access, and Consuming. A callout box pointing to Creation states:

- *Data extraction from original data sources*
- *Data transformation*
- *Metadata creation*

# Data on the Web Lifecycle



Planning → Creation → Publication

Feedback → Refining → Publication

Feedback → Consuming

*When the data (and metadata) becomes available (data catalogues, APIs)*

# Data on the Web Lifecycle



Planning → Creation → Publication

Refining → Feedback

Feedback → Consuming → Access

*When data consumers gain access to the data*

# Data on the Web Lifecycle



**Planning**

**Feedback**

**Creation**

*Data reuse (creation of new data, data visualizations and data analysis applications)*

**Consuming**

**Publication**

**Web**

**Access**

# Data on the Web Lifecycle

# Data on the Web Lifecycle



Planning

Feedback

Data maintenance (data correction and data up to date)

Refining

Consuming

Publication

Web

Access

# Players of the data on the Web ecosystem

*Several types of data sources (transactional systems, sensors, mobile devices, documents…)*

Data publisher:
publishes and shares data

Data consumer:
reuses the data and might generate new data

Source: http://ceweb.br/livros/dados-abertos-conectados/capitulo-1/

## *How to enable the data reuse?*

ceweb.br nic.br cgi.br

# How to enable the reuse of data?

*A common understanding between data publishers and data consumers becomes fundamental.*
*Without this agreement, data publishers' efforts may be incompatible with data consumers' desires.*



Consumes data

Publishes data

# Data on the Web Best Practices Working Group

The **Mission** of the Data on the Web Best Practices Working Group, part of the Data Activity, is:

1. to develop the **open data ecosystem**, <u>facilitating better communication</u> between developers and publishers;
2. to provide **guidance to publishers** that will improve consistency in the way data is managed, thus <u>promoting the re-use of data</u>;
3. to **foster trust in the data** among developers, whatever technology they choose to use, <u>increasing the potential for genuine innovation</u>.



Source: https://www.w3.org/2013/dwbp/wiki/Main_Page:

# Data on the Web Best Practices

W3C Editor's Draft 22 March 2016

**This version:**
http://w3c.github.io/dwbp/bp.html

**Latest published version:**
http://www.w3.org/TR/dwbp/

**Latest editor's draft:**
http://w3c.github.io/dwbp/bp.html

**Editors:**
Bernadette Farias Lóscio, CIn - UFPE, Braz
Caroline Burle, NIC.br, Brazil
Newton Calegari, NIC.br, Brazil

**Authors:**
Annette Greiner
Antoine Isaac
Carlos Iglesias
Carlos Laufer
Christophe Guéret
Eric G. Stephan
Eric Kauz
Ghislain A. Atemezing
Ig Ibert Bittencourt
João Paulo Almeida
Manuel Tomas Carrasco
Phil Archer
Riccardo Albertoni

## Table of Contents

*BP are designed to meet the needs of information management staff, developers, and wider groups such as scientists interested in sharing and reusing research data on the Web*



Source: http://w3c.github.io/dwbp/bp.html

ceweb.br  nic.br  cgi.br

# Data on the Web Use cases & Requirements

*scenarios of how data is commonly published on the Web and how it is used*

*cover different domains and illustrate some of the main challenges faced by data publishers and data consumers*

Source: https://www.w3.org/TR/dwbp-ucr/

# Data on the Web Challenges

- Metadata *(for humans & machines)*

- Data Licenses *(how to permite & restrict access?)*

- Data Provenance & Quality *(how to add trust?)*

- Data Versioning (*tracking dataset versions)*

- Data Identifiers *(identifying datasets and distributions)*

- Data Formats *(which data formats to use?)*

*The openness and flexibility of the Web creates new challenges for data publishers and data consumers*

# Data on the Web Challenges

- Data Vocabularies *(how to promote interoperability?)*

- Sensitive Data *(Privacy & Security)*

- Data Access *(access options)*

- Feedback *(how to engage users?)*

- Data Enrichment *(adding value to data)*

Best Practice 1: Provide metadata

Best Practice 2: Provide descriptive metadata

Best Practice 3: Provide locale parameters metadata

Best Practice 4: Provide structural metadata

Best Practice 5: Provide data license information

Best Practice 6: Provide data provenance information

Best Practice 7: Provide data quality information

Best Practice 8: Provide a version indicator

Best Practice 9: Provide version history

Best Practice 10: Use persistent URIs as identifiers of datasets

Best Practice 11: Use persistent URIs as identifiers within datasets

Best Practice 12: Assign URIs to dataset versions and series

Best Practice 13: Use machine-readable standardized data formats

Best Practice 14: Provide data in multiple formats

Best Practice 15: Use standardized terms

Best Practice 16: Reuse vocabularies

Best Practice 17: Choose the right formalization level

Best Practice 18: Provide data unavailability reference

Best Practice 19: Provide bulk download

Best Practice 20: Provide Subsets for Large Datasets

Best Practice 21: Use content negotiation for serving data available in multiple formats

Best Practice 22: Provide real-time access

**Evidence**

**Relevant requirements:** R-ProvAvailable, R-MetadataAvailable

Best Practice 26: Provide complete documentation for your API

Best Practice 27: Avoid Breaking Changes to Your API

Best Practice 28: Assess dataset coverage

Best Practice 29: Use a trusted serialisation format for preserved data dumps

Best Practice 30: Update the status of identifiers

Best Practice 31: Gather feedback from data consumers

Best Practice 32: Make feedback available

Best Practice 33: Enrich data by generating new data

Best Practice 34: Provide Complementary Presentations

Best Practice 35: Provide Feedback to the Original Publisher

Best Practice 36: Follow Licensing Terms

Best Practice 37: Cite the Original Publication

# Data on the Web Best Practices

# DWBP Benefits

*Each benefit represents an improvement in the way how datasets are available on the Web*



## Reuse

BP: Provide data license information
BP: Provide versioning information
BP: Provide version history
BP: Use non-proprietary data formats
BP: Provide data in multiple formats
BP: Use a trusted serialization format for preserved data dumps
BP: Enrich data by generating new metadata
BP: Provide data provenance information
BP: Provide data quality information
BP: Use persistent URIs as identifiers

## Trustworthy

BP: Assess dataset coverage
BP: Assign URIs to dataset versions and series
BP: Provide data up to date
BP: Update the status of identifiers
BP: Gather feedback from data consumers
BP: Provide information about feedback
BP: Provide data provenance information
BP: Provide data quality information

## Comprehension

BP: Provide metadata
BP: Provide locale parameters metadata
BP: Provide structural metadata
BP: Provide descriptive metadata

## Accessibility

BP: Provide bulk download
BP: Follow REST principles when designing APIs
BP: Provide real-time access
BP: Maintain separate versions for a data API
BP: Assess dataset coverage

## Linkability

BP: Use persistent URIs as identifiers
BP: Assign URIs to dataset versions and series

## Discoverability

BP: Provide descriptive metadata
BP: Use persistent URIs as identifiers
BP: Assign URIs to dataset versions and series

## Processibility

BP: Use machine-readable standardized data formats
BP: Enrich data by generating new metadata

## Interoperability

BP: Use standardized terms
BP: Re-use vocabularies

# BP Benefits

- Comprehension: humans will have a better understanding about the data structure, the data meaning, the metadata and the nature of the dataset.
- Processability: machines will be able to automatically process and manipulate the data within a dataset.
- Discoverability machines will be able to automatically discover a dataset or data within a dataset.
- Reuse: the chances of dataset reuse by different groups of data consumers will increase.

## Best Practice 10: Use persistent URIs as identifiers of datasets

*Datasets must be identified by a persistent URI.*

### Why

Adopting a common identification syste~~~
by any stakeholder in a reliable way. The
and reuse.

Developers may build URIs into their co~~~
dereference to the same resource over t~~~

### Intended Outcome

Datasets or information about datasets ~~~
status, availability or format of the data.

### Possible Approach to Implementation

To be persistent, URIs must be designe~~~
creating a Web site designed for human~~~
topic, see, for example, the European C~~~
to many other resources.

Where a data publisher is unable or unw~~~
native approach is to use a redirection s~~~
These provide persistent URIs that can ~~~
ephemeral. The software behind such se~~~
aged locally if required.

Digital Object Identifiers (DOIs) offer a similar alternative. These identifiers are defined independently of
any Web technology but can be appended to a 'URI stub.' DOIs are an important part of the digital infra-
structure for research data and and libraries.

# BP Benefits

- Linkability: it will be possible to create links between data resources (datasets and data items).
- Interoperability: it will be easier to reach consensus among data publishers and consumers.
- Discoverability machines will be able to automatically discover a dataset or data within a dataset.
- Reuse: the chances of dataset reuse by different groups of data consumers will increase.

# Data on the Web Best Practices

1. They are still being developed
2. Publication of the next draft as Candidate Recommendation - April 2016
3. Publication as recommendation - July 2016
4. Feedback is welcome! :)

# [https://www.w3.org/TR/dwbp/](https://www.w3.org/TR/dwbp/)

# Obrigada!

## www.ceweb.br

@ cburle@nic.br     Ⓣ @carolburle

@ bfl@cin.ufpe.br    Ⓣ @bernafarias

**April 12th, 2016**

nic.br   cgi.br

www.nic.br | www.cgi.br