First they have to **find** it:
Getting Government Data Discovered and Used

John S. Erickson, Ph.D.
Tetherless World Constellation
Rensselaer Polytechnic Institute
Troy, New York, USA

Twitter: **@olyerickson**   **#TWCRPI**
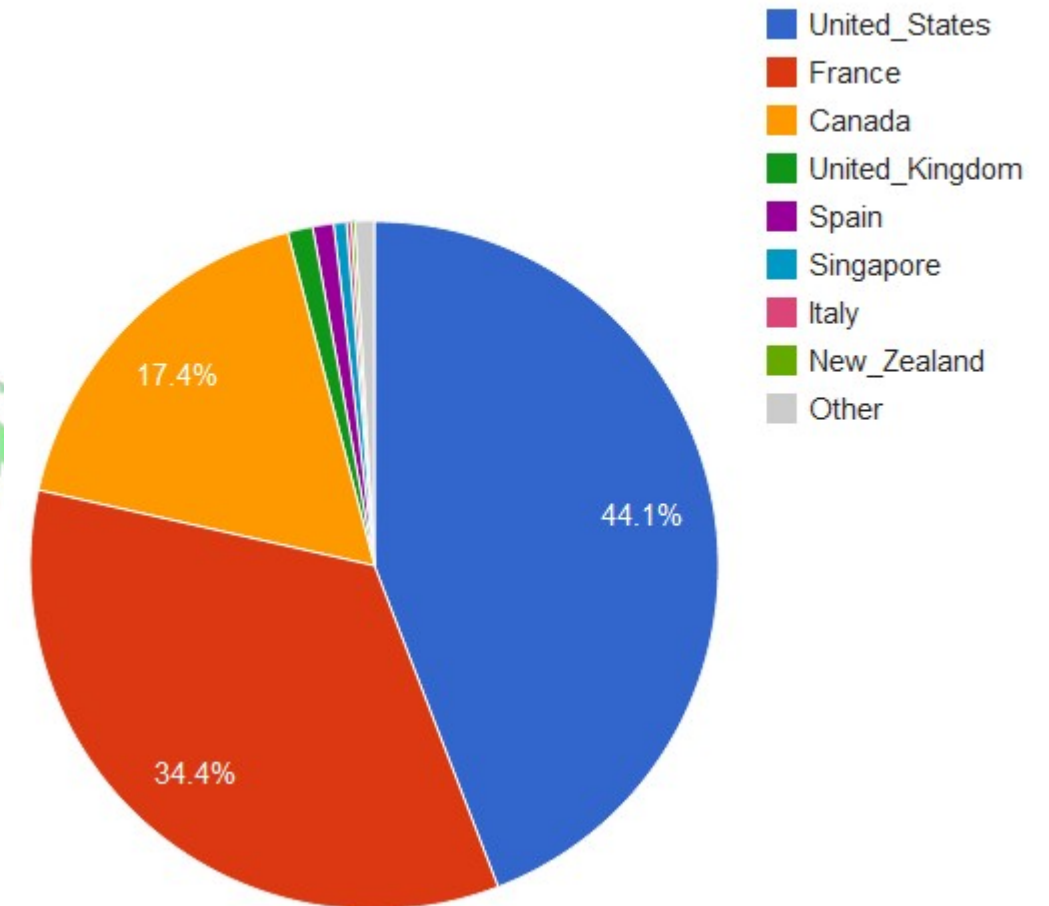
Starting with efforts in the US and UK, governments around the world have recognized the need to publish their critical data

Percent of total collection (from 1M+ datasets)



Legend:
- United_States
- France
- Canada
- United_Kingdom
- Spain
- Singapore
- Italy
- New_Zealand
- Other

Pie chart values:
- 44.1%
- 34.4%
- 17.4%

Map scale: 5 to 755,473

2

- Government data initiatives have taken many forms

- GovData portals are widely varied in how they help users **discover** and **use relevant datasets**
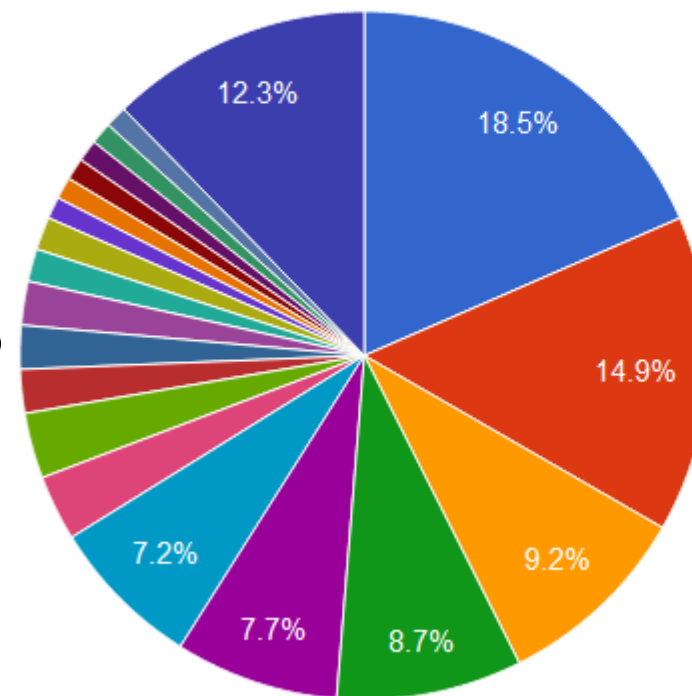
Percent of total catalogs
(from 192 catalogs)



- United_States — 18.5%
- Canada — 14.9%
- United_Kingdom — 9.2%
- Spain — 8.7%
- Italy — 7.7%
- France — 7.2%
- Austria
- Germany
- Brazil
- Australia
- Ireland
- International_organization
- Netherlands
- Belgium
- New_Zealand
- Chile
- Finland
- Slovakia
- Argentina
- The Other — 12.3%

page | discussion | view source | history

**W3C**

## Data Catalog Vocabulary

The **Data Catalog Vocabulary** (DCAT) is an RDF Schema vocabulary for metadata ab... incubated in the eGovernm...

**schema.org**

Navigation

Main page
Charter
Tracker
Email Archive
Recent changes
Help

**Contents** [hide]
1 Documents and Delivera...
2 Issue tracking
3 People
4 Goals and scope
5 Other resources

earch

Go | Search

## Organization of Schemas

The schemas are a set of 'types', each associated with a set of properties. The types are arrange...

## Documents and

Data Catalog Schema and Protocol v0.1

## Catalog Access, Federation and Harvesting Mechanism

**Status: early draft**

This portion of the specification details a protocol for accessing catalog metadata and supporting automated harvesting and federation.

*This specification is at a very early stage and is intended as a basis for discussion rather than a finished document.*

/Series

## API

A catalog MUST provide the following API. The API base location is specified by the following meta tag in the site home page:

```
<meta content="data-catalog-api" value="http://my-data-catalog.org/api" />
```
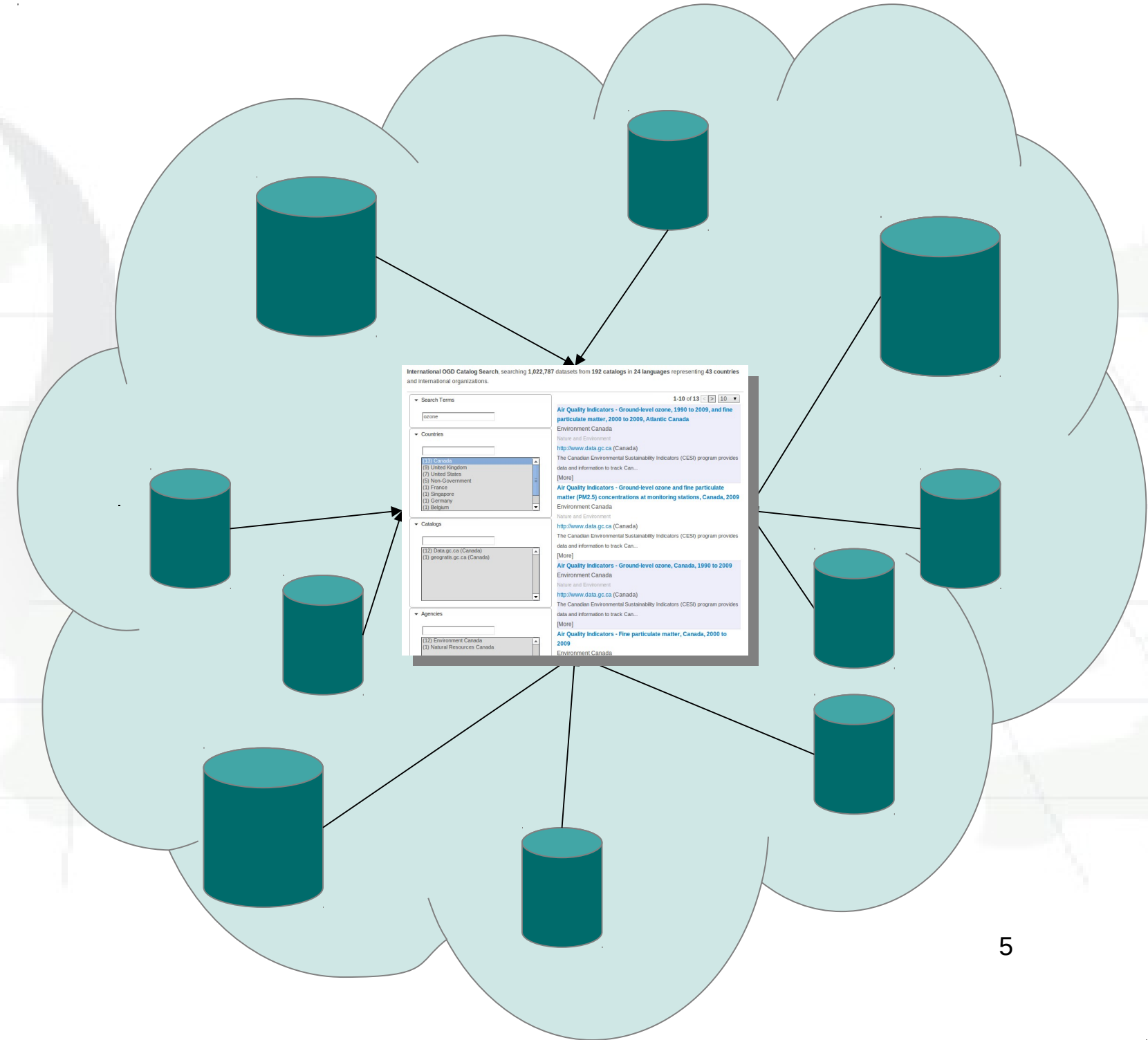
Relative to this base URL there are the following endpoints:

Stakeholders have seen the need for **Federated discovery** across catalogs, **especially** from within major search engines including **Bing, Google, Yahoo!** and **Yandex**
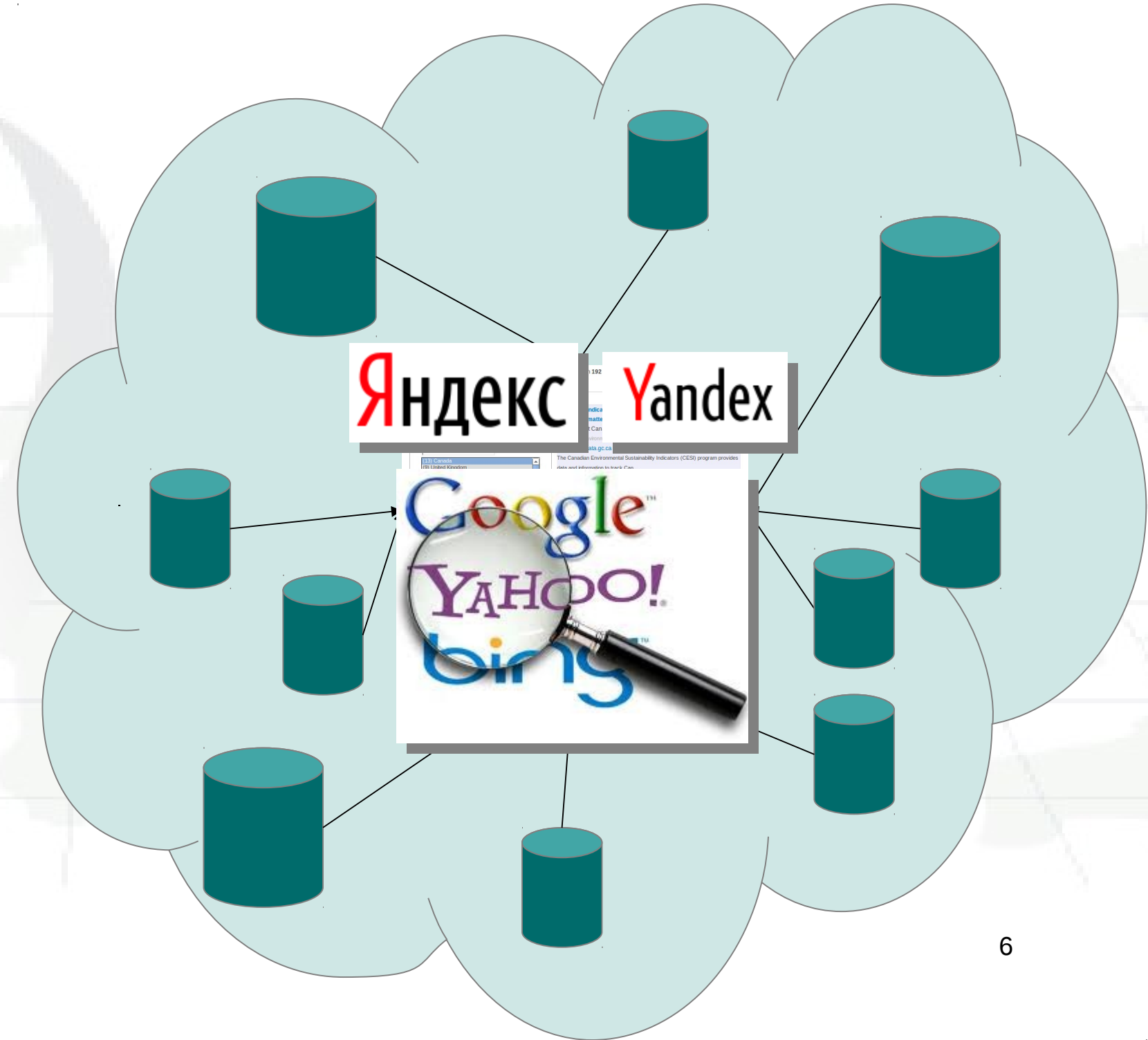
4

- Publishing open government data as Linked Data is **not enough**

- For OGD to be useful, datasets must be published using metadata, markup standards and presentation that aid **discovery** and **use**

- Publishing open government data as Linked Data is **not enough**

- For OGD to be useful, datasets must be published using metadata, markup standards and presentation that aid **discovery** and **use**

International OGD Catalog Search, searching **1,022,787** datasets from **192 catalogs** in **24 languages** representing **43 countries** and international organizations.

▼ Search Terms

ozone

1-10 of 13 < > 10 ▼

**Air Quality Indicators - Ground-level ozone, 1990 to 2009, and fine particulate matter, 2000 to 2009, Atlantic Canada**

Environment Canada

Nature and Environment

http://www.data.gc.ca (Canada)

The Canadian Environmental Sustainability Indicators (CESI) program provides data and information to track Can...

[More]

▼ Countries

(13) Canada
(9) United Kingdom
(7) United States
(5) Non-Government
(1) France
(1) Singapore
(1) Germany
(1) Belgium

▼ Catalogs

(12) Data.gc.ca (Canada)
(1) geogratis.gc.ca (Canada)

▼ Agencies

(12) Environment Canada
(1) Natural Resources Canada

Thing > CreativeWork > Dataset
A structured body of information describing some topic(s) of interest.

| Property | Expected Type | Description |
|---|---|---|
| **Properties from Thing** | | |
| description | Text | A short description of the item. |
| image | URL | URL of an image of the item. |
| name | Text | The name of the item. |
| url | URL | URL of the item. |
| **Properties from CreativeWork** | | |
| about | Thing | The subject matter of the content. |
| accountablePerson | Person | Specifies the Person that is legally accountable for the CreativeWork. |
| aggregateRating | AggregateRating | The overall rating, based on a collection of reviews or ratings, of the item. |
| alternativeHeadline | Text | A secondary title of the CreativeWork. |
| associatedMedia | MediaObject | The media objects that encode this creative work. This property is a synonym for encodings. |
| audio | AudioObject | An embedded audio object. |
| author | Person or Organization | The author of this content. Please note that author is special in that HTML 5 provides a special mechanism for indicating authorship via the rel tag. That is equivalent to this and may be used interchangably. |
| awards | Text | Awards won by this person or for this creative work. |
| comment | UserComments | Comments, typically from users, on this CreativeWork. |
| contentLocation | Place | The location of the content. |
| contentRating | Text | Official rating of a piece of content—for example,'MPAA PG-13'. |

Recent work at TWC RPI demonstrates the value of applying emerging **standards** for **uniformly describing** government datasets and catalogs

7

TWC's **IOGDS** application is an aggregated catalog of more than 1M datasets from over 192 dataset catalogs from governments at every level around the world

See: http://logd.tw.rpi.edu

W3C Editor's Draft

**W3C**

Data Catalog Vocabulary (DCAT)

W3C Editor's Draft 11 July 2012

**This version:**
http://dvcs.w3.org/hg/gld/raw-file/default/dcat/index.html
**Latest published version:**
http://www.w3.org/TR/vocab-dcat/
**Latest editor's draft:**
http://dvcs.w3.org/hg/gld/raw-file/default/dcat/index.html
**Previous version:**
none
**Editors:**
Fadi Maali, DERI, NUIG
John Erickson, Tetherless World Constellation (RPI)
Phil Archer, W3C/ERCIM

Copyright © 2012 W3C® (MIT, ERCIM, Keio), All Rights Reserved. W3C liability, trademark and document use rules apply.

**Abstract**

DCAT is an RDF vocabulary designed to facilitate interoperability between data catalogs published on the Web. This document defines the schema and provides examples for its use.

By using DCAT to describe datasets in data catalogs, publishers increase discoverability and enable applications easily to consume metadata from multiple catalogs. It further enables decentralized publishing of catalogs and facilitates federated dataset search across sites. Aggregated DCAT metadata can serve as a manifest file to facilitate digital preservation.
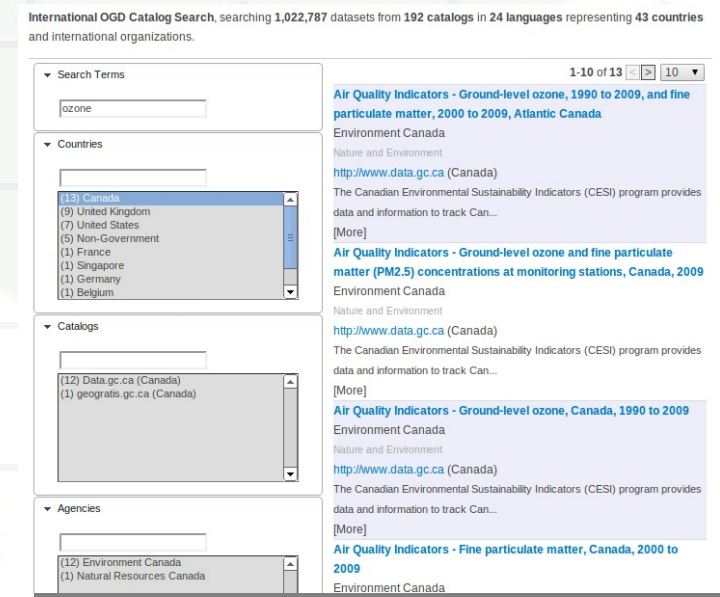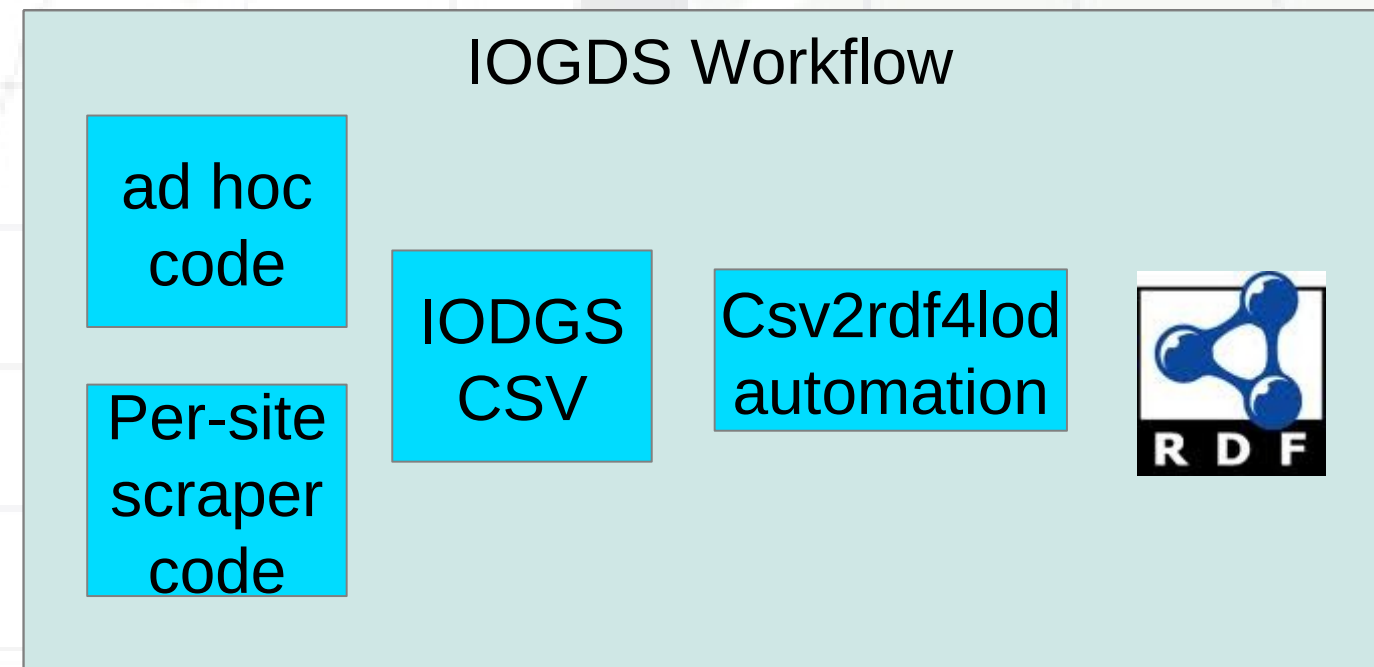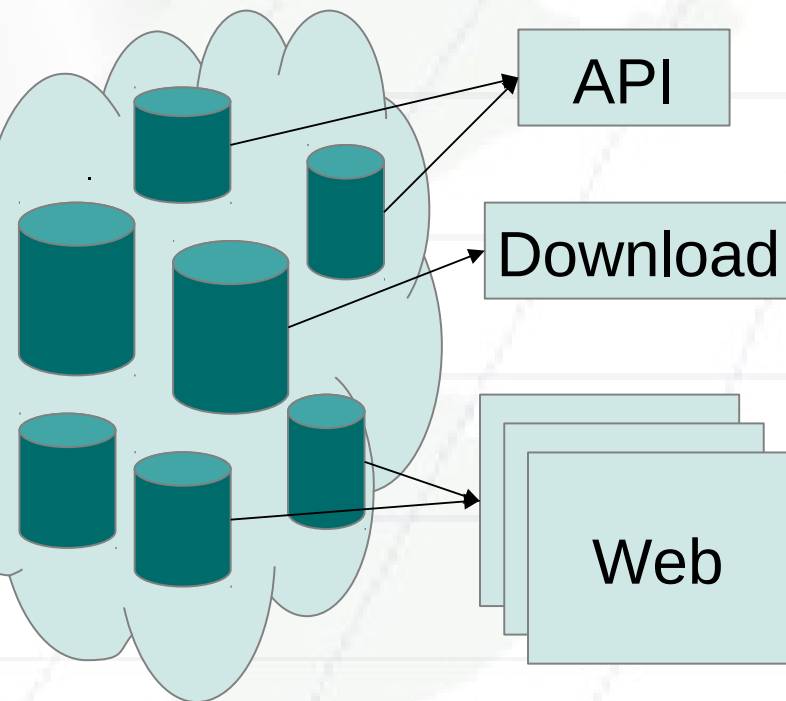
- Anticipates W3C DCAT RDF vocabulary
- Demos what a comprehensive **federated catalog** based on DCAT and aggregation API might look like

9

IOGDS is a multi-year effort based on downloading, scraping or accessing APIs, converting metadata to a proto-DCAT model, and publishing via endpoint and download

Catalogs

API

Download

Web

IOGDS Workflow

ad hoc code

Per-site scraper code

IODGS CSV

Csv2rdf4lod automation

R D F

International OGD Catalog Search, searching 1,022,787 datasets from 192 catalogs in 24 languages representing 43 countries and international organizations.

See: http://logd.tw.rpi.edu

**Seismic Hazard Zones**

The dataset represents the Liquefaction and Landslide Zones as determined bt the California Dept. of Conservation Division of Mines and Geology. Liquefaction is the transformation of a confined layer of sandy or silty water-satuated material into a liquid -like state because of earthquake shaking. San Francisco Building Code Section1804.5 requires a geotechnical investigation in seismic hazard zones.
*Country:* United States
*Publisher:* Department of Technology
*Categories:* layers , geography, maps , gis

**Street Sweeper Schedule and Route**

Street sweeper schedule that includes: when and where.
*Country:* United States
*Publisher:* Department of Public Works
*Categories:* department of public works , street sweeper

**Street Centerlines**

View of Street Centerlines excluding Paper streets, unpaved rights-of-way and psuedo streets.
*Country:* United States
*Publisher:* Department of Public Works
*Categories:* layers , geography, maps , gis

**USPS Post Offices**

USPS Post Offices in San Francisco.
*Country:* United States
*Publisher:* Department of Technology
*Categories:* layers , geography, maps , gis

**SF Shoreline**

San Francisco mainland shoreline and in the south, the county line.
*Country:* United States
*Publisher:* Department of Public Works
*Categories:* layers , geography, maps , gis

TWC RPI has published dataset listings based on IOGDS using emerging microdata standards, esp. **schema.org** model endorsed by **Bing, Google, Yahoo!, Yandex...**

11

# Schema.org **datasets** extension

- TWC RPI's schema.org **dataset extension** will enable government dataset catalogs to more easily be parsed and indexed by the major search engines...

- ...which will help users find relevant datasets!

- TWC's dataset extension entered public discussion June 2012

The schema.org **datasets** extension enables **relevant datasets** to be more easily **discovered** by a range of stakeholders including researchers, data journalists, bloggers and developers

13

**Seismic Hazard Zones**

The dataset represents the Liquefaction and Landslide Zones as determined bt the California Dept. of Conservation Division of Mines and Geology. Liquefaction is the transformation of a confined layer of sandy or silty water-satuated material into a liquid -like state because of earthquake shaking. San Francisco Building Code Section1804.5 requires a geotechnical investigation in seismic hazard zones.
*Country:* United States
*Publisher:* Department of Techtology
*Categories:* layers , geography, maps , gis

**Street Sweeper Schedule and Route**

Street sweeper schedule that includes:
*Country:* United States
*Publisher:* Department of Public Works
*Categories:* department of public works

**Street Centerlines**

View of Street Centerlines excluding Pa
*Country:* United States
*Publisher:* Department of Public Works
*Categories:* layers , geography, maps ,

**USPS Post Offices**

USPS Post Offices in San Francisco.
*Country:* United States
*Publisher:* Department of Technology
*Categories:* layers , geography, maps ,

**SF Shoreline**

San Francisco mainland shoreline and
*Country:* United States
*Publisher:* Department of Public Works
*Categories:* layers , geography, maps ,

**DATA.GOV / SEMANTIC**
EMPOWERING PEOPLE

Home | Apps | Blogs | Forums | Search

Data.gov » All Communities » Semantic Web » Blogs

## SUPPORT FOR SCHEMA.ORG AT DATA.GOV

Posted on Sat, 2012-07-07 19:48 by Chris Musialek

SHARE

0
0

We've been watching the schema.org datasets schema space for a while now, as Data.gov is very interested in adding schema.org support for our listing of over 450,000 datasets. We think this will help the major search engines create better relevance rankings of Federal government data, where many searches begin.

We wanted to come out publicly saying that we've reviewed the current datasets schema **proposal** in draft, and we are comfortable with the current state of things. There is definitely work still left to do, but there seems to be pretty solid agreement on everything but the details, which seem very resolvable. At this point, if the group would solidify on the dataset proposal, then Data.gov would support and use it.

We're really excited to see this schema move in the direction of official addition to schema.org. We really hope to see it be included in a schema.org release soon.

"...we've reviewed the current datasets schema proposal in draft, and we are comfortable with the current state of things...

"...At this point, if the group would solidify on the dataset proposal, then **Data.gov would support and use it.**

---Chris Musialek

**DATA.GOV**
EMPOWERING PEOPLE

# CKAN Data Catalog Scheme & Protocol

- API-based catalog federation is also possible

- ckan announced DCAT-based query/federation API

- enables OAI-PMH-like harvesting and more

Data Catalog Schema and Protocol v0.1

## Catalog Access, Federation and Harvesting Mechanism

**Status: early draft**

This portion of the specification details a protocol for accessing catalog metadata and supporting automated harvesting and federation.

*This specification is at a very early stage and is intended as a basis for discussion rather than a finished document.*

## API

A catalog MUST provide the following API. The API base location is specified by the following meta tag in the site home page:

```
<meta content="data-catalog-api" value="http://my-data-catalog.org/api" />
```

Relative to this base URL there are the following endpoints:

```
/changes.json # changes API
/dataset/{id}.json # dataset API
```

- Geo-based discovery:
  *What data is available by geo-selection?*

- Provenance-based discovery:
  *How do I get the data that someone else used? "Get the Data"*

- Community/social-based discovery:
  *Dude, check out this data! (Linked Data perfect for this...*

# Other Thoughts...

- Geo-based discovery:
  *What data is available by geo-selection?*

  **DATA.gov** *Geo Viewer*

# Other Thoughts...

- Community/social-based discovery:
*Dude, check out this data!*

**OPENEI.org**

Choose your own medicine...
but **do** expose your **metadata**
and **get** your catalogs **discovered**!