

Exploiting Open Data

David Tarrant
@davetaz

Toolkit for today

1. Google Chrome Web Browser
chrome.google.com
2. Open/Google Refine
openrefine.org
3. A Google Account
4. The Postman extension to Google Chrome (pm only)

Aims (am)

- Discovering, filtering, combining and analyzing data
- Reconciliation, enriching and linking data
- Small graphs from “Big Data”

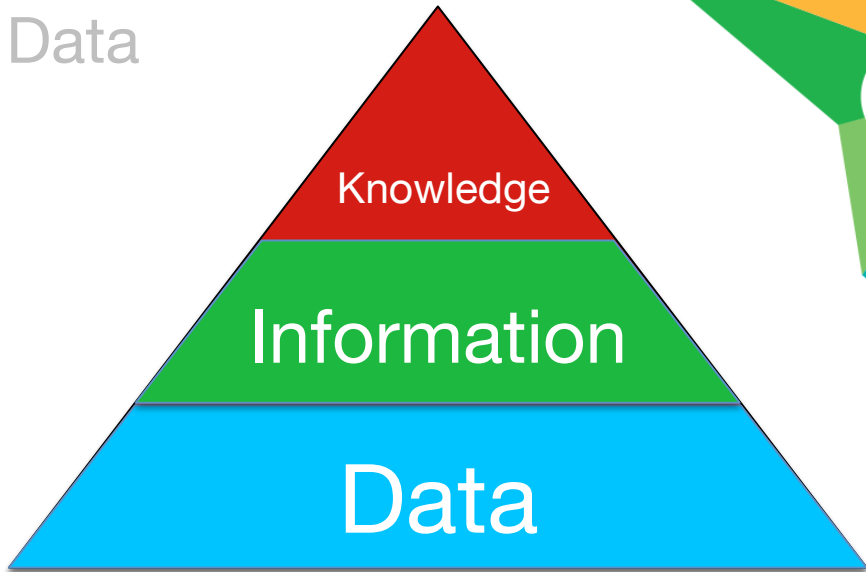


Recap

What is Open Data?



Data



odi

Definition of Open (OKF)



A piece of data or content is open if **anyone** is **free to use, reuse, and redistribute** it — subject only, at most, to the requirement to attribute and/or share-alike.

Trends in publishing open data

Suppliers examples



**Private sector
open datasets**

“While there is no central figure on the number of public sector information datasets currently being made available, a review of selected data portals suggests the number could exceed 37,500 from over 750 different publishers with over 2.5 million downloads by October 2012” Deloitte 2012

Note

Open data examples:



Non open data examples:



Suppliers

bp

BP: GBP 506.2 ⬆ (-0.8) USD 51.1 ⬆ (-0.11) *

BP Global | BP Worldwide

About BP | Products and Services | Sustainability | Investors | Press | Careers | Gulf of Mexico restoration

About BP > Energy economics > Statistical Review of World Energy 2013 > Downloads

Statistical Review downloads

Use this section to download the Excel workbook of historical data, the printed edition and the launch presentation speech and slides

Statistical Review 2013

Download the printed edition of the Statistical Review of World Energy 2013, the supporting PowerPoint Excel workbook of historical statistical data from 1965-2012

- Statistical Review of World Energy 2013 (pdf, 9.6MB)
- Statistical Review 2013 workbook (xlsx, 1.5MB)**
- Statistical Review slidepack (pptx, 17MB)
- Statistical Review 2013 Speech (pdf, 663.0KB)
- Statistical Review speech slide pack (pptx, 2.9MB)

Energy Outlook

OPEN DATA



BP: You may not frame this site nor link to a page other than the home page without our express permission.
To find this Google “bp statistical review”

Suppliers (ish)



Bigger Picture

Approach

Our contribution

Responsible business

Inspiring action

News, views & videos

How we are doing

How we report

Data commentary

Financial data

Creative industries data

UK economy data

Customers data

People data

Business partners data

Environmental impact data

Community data

Sport data

Environmental impact data



At Sky we understand that to build a sustainable business we need to minimise our environmental impact from our operations, create more sustainable products and services and work with our business partners and customers to inspire them to take action to protect the environment too.

Our progress so far

Last year we reported good progress towards our ten environment targets that were set in 2009 so we set new, more ambitious targets to 2020. This year we have continued to reduce our emissions when normalised against turnover. We have achieved this by investing in energy efficiency measures, working with our people to reduce the energy they use day to day and through our on-site renewable energy sources. Detailed commentary about this data is provided below.

Summary performance against environment targets

	Target	2009/10	2010/11	2011/12	2012/13
Reduction in gross CO ₂ e emissions relative to revenue(%) ¹	-50	-8	-21	-29	-33
Energy obtained from owned or controlled renewable energy at Sky-owned sites(%) ²	20	*	*	2	
Increase in fleet fuel efficiency(%) ³	-15	*	*	*	-5
Reduction in CO ₂ e					

OPEN DATA

Suppliers (ish)

Product Environment Reports



Date introduced

September 10, 2013

Environmental Status Report



iPhone 5s is designed with the following features to reduce environmental impact:

- Arsenic-free display glass
- Mercury-free LED-backlit display
- Brominated flame retardant-free
- PVC-free

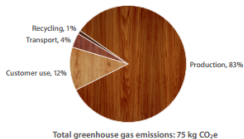
Apple and the Environment

Apple believes that improving the environmental performance of our business starts with our products. The careful environmental management of our products throughout their life cycles includes controlling the quantity and types of materials used in their manufacture, improving their energy efficiency, and designing them for better recyclability. The information below details the environmental performance of iPhone 5s as it relates to climate change, energy efficiency, material efficiency, and restricted substances.¹

Climate Change

Greenhouse gas emissions have an impact on the planet's balance of land, ocean, and air temperatures. Most of Apple's corporate greenhouse gas emissions come from the production, transport, use, and recycling of its products. Apple seeks to minimize greenhouse gas emissions by setting stringent design-related goals for material and energy efficiency. The chart below provides the estimated greenhouse gas emissions for iPhone 5s over its life cycle.²

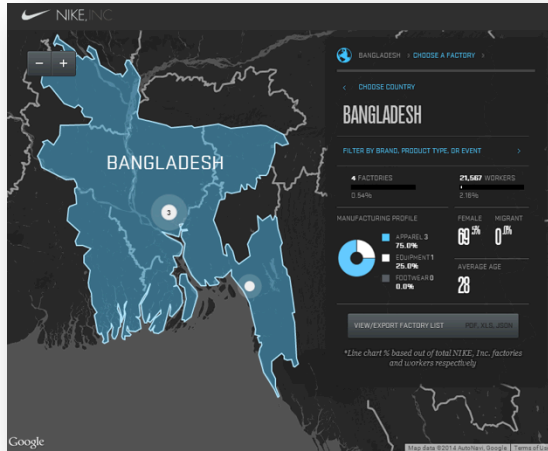
Greenhouse Gas Emissions for iPhone 5s



ODI

OPEN DATA

Suppliers



<http://manufacturingmap.nikeinc.com/#>

You agree not to change or delete any ownership notices from materials downloaded or printed from the Platform. You agree not to modify, copy, translate, broadcast, perform, display, distribute, frame, reproduce, republish, download, display, post, transmit or sell any Intellectual Property or Content appearing on the Platform

Government

Statistics

Geography

Statistics by topic - Google Chrome

Home > Statistics>

Población, Hogar

- Summary table
- Population
 - Size and growth
 - Distribution by age and sex
- Birth and fertility
- Marriage
- Migration
- Mortality
- Homes
- Dwelling

Format:

Excel 5.0 (. Xls) ▼
Excel 5.0 (. Xls)
Word (. Doc)

Export

Population Density

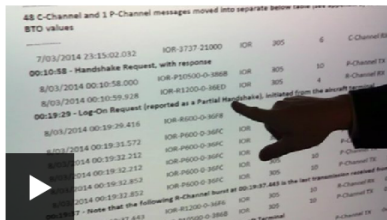
Federal District	474 860	541 516	720 753	906 063	1229576	1757530	3050442	4870876	6874165	8831079	8235744	8489007	8605239	87
Durango	296 979	370 294	483 175	336 766	404 364	483 829	629 874	760 836	939 208	1182320	1349378	1431748	1448661	15
	370	3006110	3982593	4406568	4663032	48								
	360	2109513	2620637	2916567	3079649	31								
	345	1547493	1888366	2112473	2235591	23								
	586	4371998	5302689	5991176	6322002	67								
	185	7564335	9815795	11707964	13096686	140								
	226	2868824	3548199	3870604	3985667	39								
	119	947 089	1195059	1442662	1555296	16								
	031	726 120	824 643	896 702	920 185	9								
	389	2513044	3098736	3550114	3834141	41								
	424	2369076	3019560	3228895	3438765	35								
Puebla	992 426	1021133	1101600	1024955	1150425	1294620	1625830	1973837	2508226	3347685	4126101	4624365	5076686	53
Querétaro	232 305	232 389	244 663	220 231	234 058	244 737	286 238	355 045	485 523	739 605	1051235	1250476	1404306	15
Quintana Roo	NA	NA	9109	10,966	10,620	18,752	26,967	50,169	88,150	225 985	493 277	703 536	874 963	11
San Luis Potosí	571 420	575 432	627 800	445 681	579 831	678 779	856 066	1048297	1281996	1673893	2003187	2200763	2299360	24
Sinaloa	261 050	296 701	323 642	341 265	395 618	492 821	635 681	838 404	1266528	1849879	2204054	2425675	2536844	26
Sonora	192 721	221 682	265 383	275 127	316 271	364 176	510 607	783 3						
Tabasco	134 956	159 834	187 574	210 437	224 023	285 630	362 716	496 3						
Tamaulipas	209 106	218 948	249 641	286 904	344 039	458 832	718 167	10241						
Tlaxcala	168 358	172 315	184 171	178 570	205 458	224 063	284 551	346 6						

OPEN DATA

Government / Private (ish)



Flight MH370: Malaysia releases raw satellite data



The BBC's Richard Westcott visited Inmarsat's headquarters to find out what the data tells us about MH370's fate

The Malaysian government has released the raw data used to determine that the missing Malaysia Airlines flight MH370 crashed into the southern Indian Ocean.

The data was first released to relatives of passengers, who have been asking for greater transparency, before copies were also provided to media.

The document released on Tuesday comprises 47 pages of data, plus notes, from British firm Inmarsat.

Time	Channel Name	Ocean Region	GES ID (octal)	Channel Unit ID	Channel Type	SU Type	Burst Frequency Offset (Hz) BFO	Burst Timing Offset (microseconds) BTO
03/03/2014 16:00:13.406	IOR-A1200-0-3603	IOR	305	8	R-Channel RX	0x15 - Log-on/Log-off Acknowledge		
03/03/2014 16:00:13.906	IOR-P10500-0-3859	IOR	305	10	P-Channel TX	0x15 - Log-on/Log-off Acknowledge	103	14820
03/03/2014 16:00:17.430	IOR-A1200-0-3603	IOR	305	8	R-Channel RX	Eleven Octet User Data		
03/03/2014 16:00:17.906	IOR-A1200-0-3603	IOR	305	8	R-Channel RX	Eleven Octet User Data		
03/03/2014 16:00:18.406	IOR-A1200-0-3603	IOR	305	8	R-Channel RX	Eight Octet User Data		14740
03/03/2014 16:00:18.906	IOR-P10500-0-3859	IOR	305	10	P-Channel TX	0x12 - Acknowledge User Data	103	14780
03/03/2014 16:00:20.906	IOR-P10500-0-3859	IOR	305	10	P-Channel TX	0x71 - User Data (SU) - RLS	103	14820
03/03/2014 16:00:20.906	IOR-P10500-0-3859	IOR	305	10	P-Channel TX	Subsequent Signalling Unit		
03/03/2014 16:00:20.906	IOR-P10500-0-3859	IOR	305	10	P-Channel TX	Subsequent Signalling Unit		
03/03/2014 16:00:23.407	IOR-A1200-0-3603	IOR	305	8	R-Channel RX	Subsequent Signalling Unit		
03/03/2014 16:00:23.906	IOR-P10500-0-3859	IOR	305	10	P-Channel TX	Subsequent Signalling Unit		
03/03/2014 16:00:27.741	IOR-T1200-0-3607	IOR	305	8	T-Channel RX	Subsequent Signalling Unit	88	9820
03/03/2014 16:00:27.901	IOR-T1200-0-3607	IOR	305	8	T-Channel RX	Subsequent Signalling Unit	88	9820
03/03/2014 16:00:28.061	IOR-T1200-0-3607	IOR	305	8	T-Channel RX	Subsequent Signalling Unit	88	9820
03/03/2014 16:00:28.221	IOR-T1200-0-3607	IOR	305	8	T-Channel RX	Subsequent Signalling Unit	88	9820
03/03/2014 16:00:28.405	IOR-T1200-0-3607	IOR	305	8	T-Channel RX	Subsequent Signalling Unit	88	9820
03/03/2014 16:00:28.541	IOR-T1200-0-3607	IOR	305	8	T-Channel RX	Subsequent Signalling Unit	88	9820

Government

Land Registry Monthly Property Transaction Data

Published by Land Registry. Licensed under **OGL** Open Government Licence.

Openness rating: ★★★★★

Transaction data gives numbers of applications for first registrations, leases, transfers of part, dealings, official copies and searches lodged with Land Registry by account holders in the preceding month. The information is divided into data showing all applications lodged, transactions for value, by region and local authority district. Transactions for value include freehold and leasehold sales. The data published on this page gives you information about the number and types of applications that we have completed. The data reflects the volume of applications lodged by customer Registry account number on their application form. The data does not include applications that are not yet completed, or all the applications lodged

We strive to ensure that our data is as accurate as possible but cannot guarantee that it is free from errors: <http://www.landregistry.gov.uk/public/info/data/faqs#m4>

DATA RESOURCES (151 IN A TIME SERIES)



CSV

31/12/2011 Number and types of applications by all account customers



CSV

31/12/2011 Number of searches by all account customers



CSV

31/12/2011 Number of applications in England and Wales divided by local authority district



CSV

31/12/2011 Number of applications in England and Wales divided by region

Exercise

Analyse a number of different suppliers of data and critique the different methods through which the data is available.
What are the benefits and drawbacks of each?

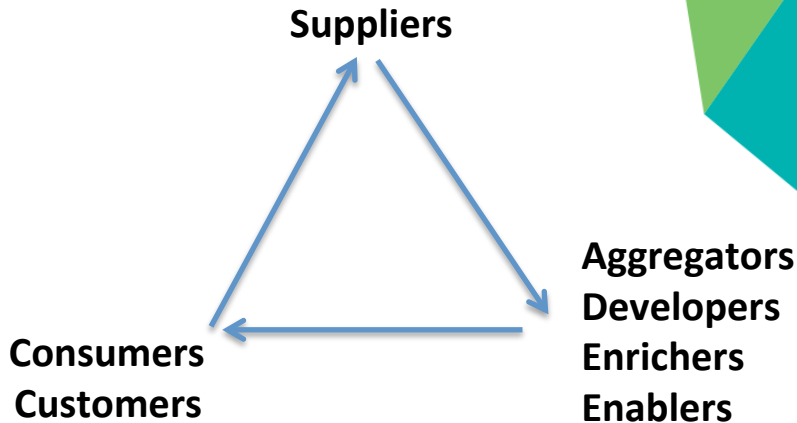


Current companies

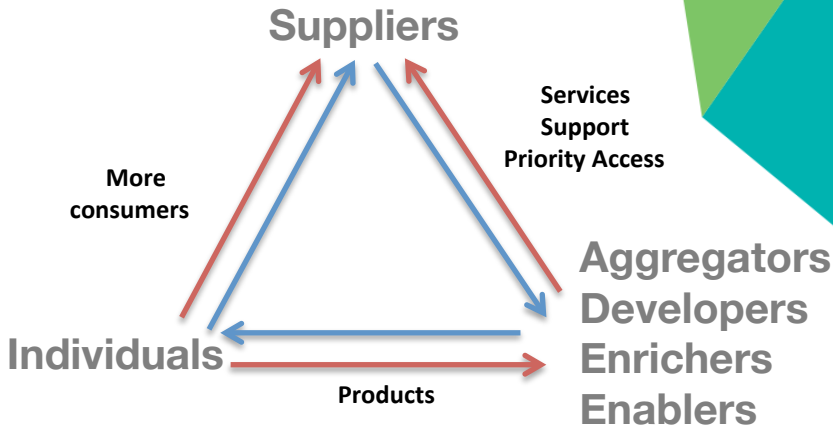
- Suppliers
- Aggregators
- Developers
- Enrichers
- Enablers



Value chain



Money go round?



Easy win: Publishing ESG (environmental, social, governance) data

More companies are deciding that the best strategy is to operate sustainably and release open data that shows it

- Sustainable practices as a sign of good corporate governance & predictor of long-term profitability
- Reduction in investment risk and helps attract new investment
- Good for branding and recruiting
- Can improve operations
- Be a step ahead of Government regulation in this area

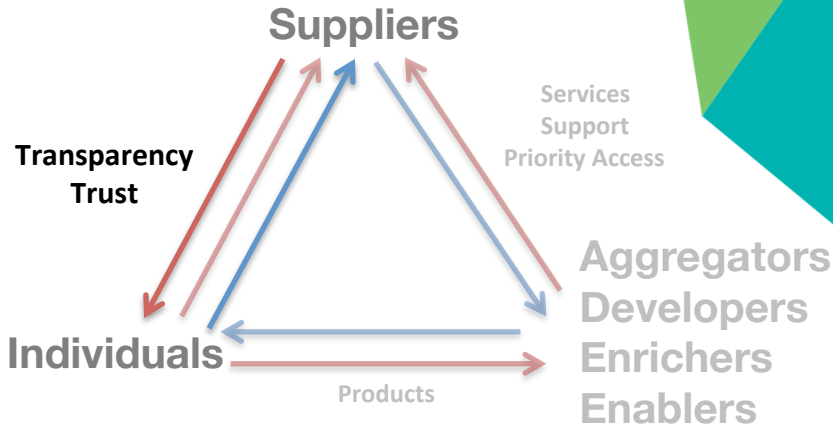
Where can businesses start **using** data and **publishing** it openly?

Increase **transparency** / **trust** / **reputation**

Operate more **efficiently** and make more informed decisions

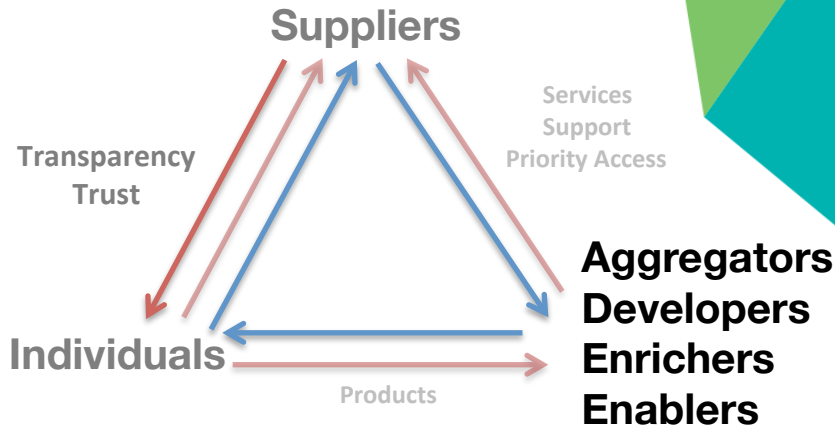
Innovate to deliver new products, services and ways of working

Money go round?



Data is the raw material of the
new industrial revolution

Money go round?



OFFICE OF FOREIGN LABOR CERTIFICATION

H-1B Visa Applications 2013

The H-1B is a non-immigrant visa in the United States under the Immigration and Nationality Act that allows U.S. employers to temporarily employ foreign workers in specialty occupations. [Full Description](#)

[+ See all nodes](#)

- UNITED STATES
- U.S. FEDERAL GOVERNMENT
- DEPARTMENT OF LABOR
- OFFICE OF FOREIGN LABOR CERTIFICATION
- H-1B VISA APPLICATIONS
- H-1B VISA APPLICATIONS 2013

H-1B VISA APPLICATIONS 2013

2,167 OF 442,277 ROWS

SHARE

EXPORT

[Add filter...](#) |

LCA Case Number	Status	Job Title	Employer Name	Employer Address	Employer City
I-200-12271-179543	DENIED	SOFTWARE ENGINE...	GOOGLE INC.	1600 AMPHITHEA...	MOUNTAIN VIEW
I-200-12251-455849	CERTIFIED	SOFTWARE ENGINE...	GOOGLE INC.	1600 AMPHITHEA...	MOUNTAIN VIEW
I-200-12265-043866	CERTIFIED	BUSINESS ANALYST	GOOGLE INC.	1600 AMPHITHEA...	MOUNTAIN VIEW
I-200-12268-519668	CERTIFIED	SOFTWARE ENGINE...	GOOGLE INC.	1600 AMPHITHEA...	MOUNTAIN VIEW
I-200-12268-506660	CERTIFIED	WEB DEVELOPER	GOOGLE INC.	1600 AMPHITHEA...	MOUNTAIN VIEW
I-200-12269-701878	CERTIFIED	SOFTWARE ENGINE...	GOOGLE INC.	1600 AMPHITHEA...	MOUNTAIN VIEW
I-200-12265-495825	CERTIFIED	INFORMATION SEC...	GOOGLE INC.	1600 AMPHITHEA...	MOUNTAIN VIEW

Exercise

- Enigma.io
- Selecting, filtering and cleaning data
- Using pivot tables and open refine.



Application Programming Interfaces (APIs)

Exercise

What is an API?
Any examples?



What is an API

Defines how one application can **consistently** interact with another.



Teaching nerds passion

```
function makeOut(passionLevel, partsOfBody) {  
  for (each partOfBody in partsOfBody) {  
    partOfBody.kiss(passionLevel);  
    lookIntoEyes();  
    sighDeeply();  
  }  
  moanDaintily();  
  sleep();  
}
```

Making your call

```
function makeOut(passionLevel, partsOfBody) {  
  for (each partOfBody in partsOfBody) {  
    partOfBody.kiss(passionLevel);  
    lookIntoEyes();  
    s makeOut(10, ["neck", "ear", "mouth"]);  
  }  
  moanDaintily();  
  sleep();  
}
```


The Webs API

```
class Photos {  
    function search(api_key, tags) {  
        if (!validKey(api_key)) {return(401);}   
        ...  
        ...  
        return results;  
    }  
}
```

Using the Webs API

[http://api.flickr.com/services/rest/?
method=flickr.photos.search&api_key=a8c42ef2ac8d8
8aed7e351f95e93160f&tags=fox](http://api.flickr.com/services/rest/?method=flickr.photos.search&api_key=a8c42ef2ac8d88aed7e351f95e93160f&tags=fox)

Evolution of Web APIs

- Long URLs are not easy to understand
- Flickr API using query parameters to define everything!

`?method=...&auth=...&query=...`

- There are better ways

Web APIs (before)

[http://api.flickr.com/services/rest/?
method=flickr.photos.comments.getList&api
_key=a8c42ef2ac8d88aed7e351f95e93160f&pho
to_id=13992401185](http://api.flickr.com/services/rest/?method=flickr.photos.comments.getList&api_key=a8c42ef2ac8d88aed7e351f95e93160f&photo_id=13992401185)

A better solution?

```
https://www.flickr.com/photos/{user}/{photo}/comments
                                     /tags
                                     /licence
```

<https://www.flickr.com/search?tags=...>

API Recap

A **promise** by one application to service another.

The underlying code can change separately to the API...

So far we have looked at READ only





Aggregator/Enricher

- Take data from hundreds of transport companies
- Bring it together
- Align it under one easy to use API

We have information on
70,597,888 companies

 SEARCH

☒ search companies ☐ search officers

Filter by jurisdiction

1,298 [Abu Dhabi \(UAE\)](#)
144,755 [Alaska \(US\)](#)
40,157 [Albania](#)
899,455 [Arizona \(US\)](#)
46,537 [Aruba](#)
165,582 [Bahamas](#)
99,185 [Bahrain](#)
88,563 [Bangladesh](#)

Just released:
OpenCorporates API v0.3

Corporate network data,
financial accounts, complex
filters, and more. [Read more](#)

Get data access to over
60 million companies

Open data	Quality data	Unique data
<ul style="list-style-type: none"> All the data on the world's largest open database of companies Available in many formats: CSV, JSON, XML, etc. API access to OpenCorporates data API access to OpenCorporates data 	<ul style="list-style-type: none"> Highly accurate, verified, and audited data Highly accurate, verified, and audited data Highly accurate, verified, and audited data Highly accurate, verified, and audited data 	<ul style="list-style-type: none"> Open, transparent, and verifiable data Open, transparent, and verifiable data Open, transparent, and verifiable data Open, transparent, and verifiable data

Announcing Open LEIs

Today, OpenCorporates announces a new sister website, [Open LEIs](#), a user-friendly interface on the emerging Global Legal Entity Identifier System. [Read more](#)

OPENLEIs
A BETA VIEW ON THE LEI SYSTEM

New! Just added: Open corporate network data

[Read more](#) about this important new feature



Exercise

Enriching data

Taking the FCO dataset and using
open corporates in Refine



Discussion

What's the biggest problem with APIs like flickr and transport API?

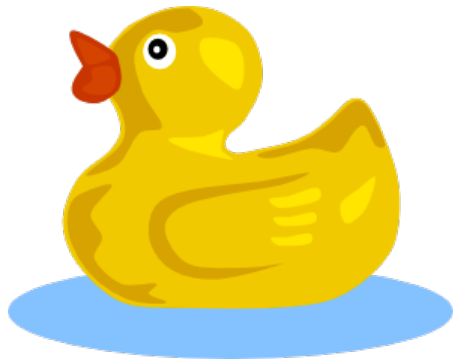


Common data retrieval methods

Duck typed data

If it looks like a duck
and quacks like a duck,
then it's probably a duck.

Basically, keep an eye out for tables,
lists and other stuff that looks like data.



Getting hold of data

1. Add random extensions (.xml, .json, .csv, .rdf etc)
2. Look for alternative links
3. Look for embedded data
4. Do some content negotiation
5. Spot the API
6. Abandon hope...(search google again)

Getting hold of data

1. Add random extensions (.xml, .json, .csv, .rdf etc)
2. Look for alternative links
3. Look for embedded data
4. Do some content negotiation
5. Spot the API
6. Abandon hope...(search google again)

1. Add random extensions

(.xml, .json, .csv, .rdf etc)

Try out a number of the links listed on the website to see if you can find data.

- When did you?
- What formats?
- Machine readable license?
- Is the URL an identifier?

2. Look for link alternates

These are in the source of the HTML page

```
<link rel="alternate" type="text/csv" href="http..."  
title="Alternate for ..."/>
```

Find these on dbpedia, CKAN, gov.uk, eprints.

3. Look for embedded data

Also in the source of the HTML page

Harder to spot as there are many different variants.

Watch out for a number of common namespace declarations in the <html> tag:

```
<html xmlns:foaf="http://xmlns.com/foaf/0.1/"  
xmlns:owl="http://www.w3.org/2002/07/owl#" xmlns:rdfs="http://  
www.w3.org/2000/01/rdf-schema#" xmlns:dc="http://purl.org/dc/  
terms/" xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-  
ns#" xmlns:dcat="http://www.w3.org/ns/dcat#" </html>
```

4. Do some content negotiation

This is the **correct answer** and we will come back to it later when we talk about the web.

Basically Tim had the answer in 1994 and we all forgot about it and invented a load of other crap instead!

Recap

- The supply of open data
- Tools for aggregating, processing and enriching data.
- Advanced data discovery
- APIs, extensions, embedded data and links.



One more thing

The promise of “Big” data.

<https://data.cityofchicago.org/Public-Safety/Crimes-2001-to-present/ijzp-q8t2>

Exercise





David Tarrant
@davetaz