# Schedule

Day 1: Publishing open data

Day 2: Business, the law and open data

Day 3: Using, enriching and visualising data

# Agenda - Today

The characteristics of data

Data discovery patterns

*** Lunch ***

Data publication platforms

*Quick big data break*

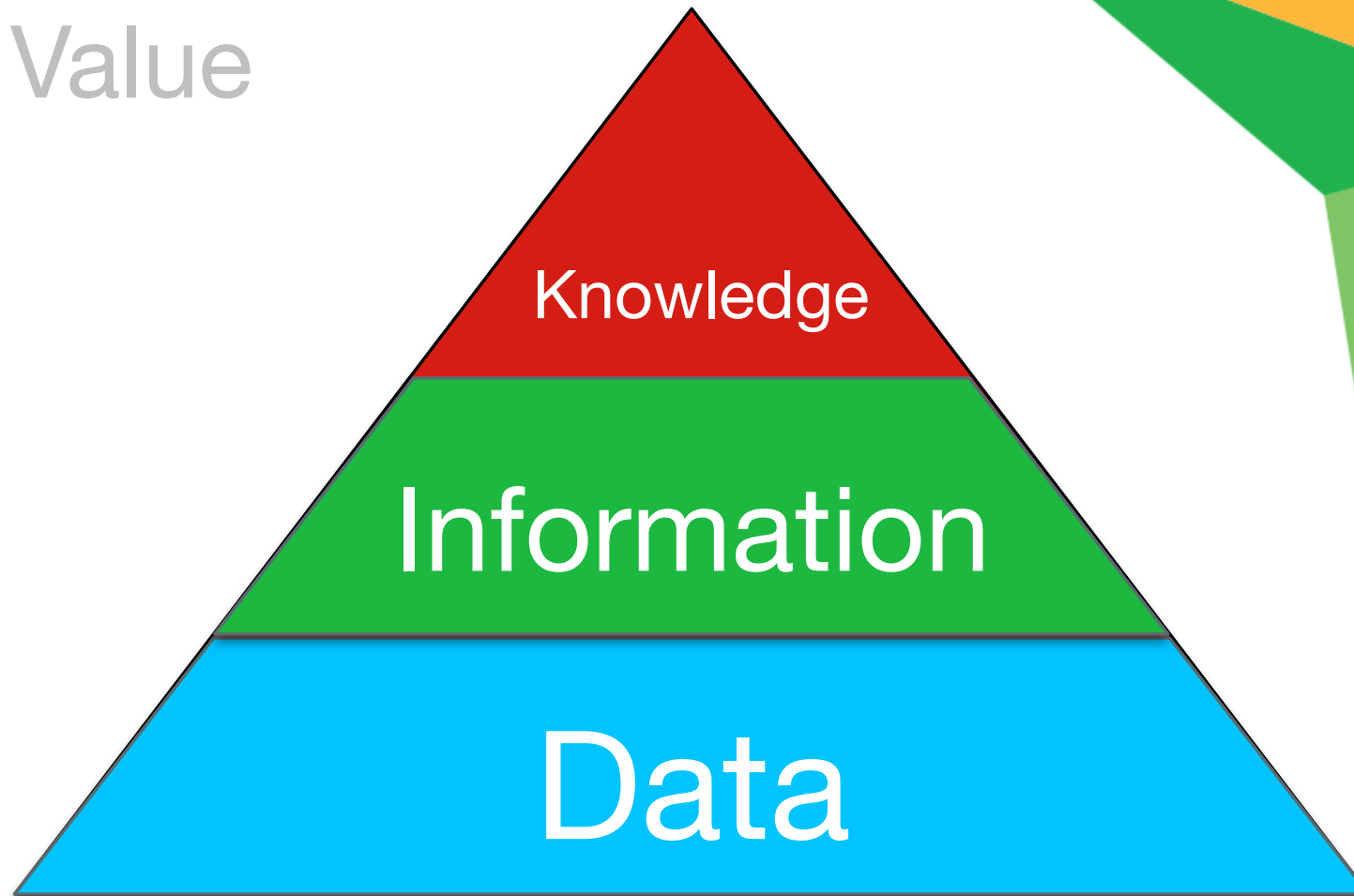Practical publication hands-on

# Recap session

Exercise

# What is Data?

Value

Knowledge

Information

Data

Exercise

What is Open Data?
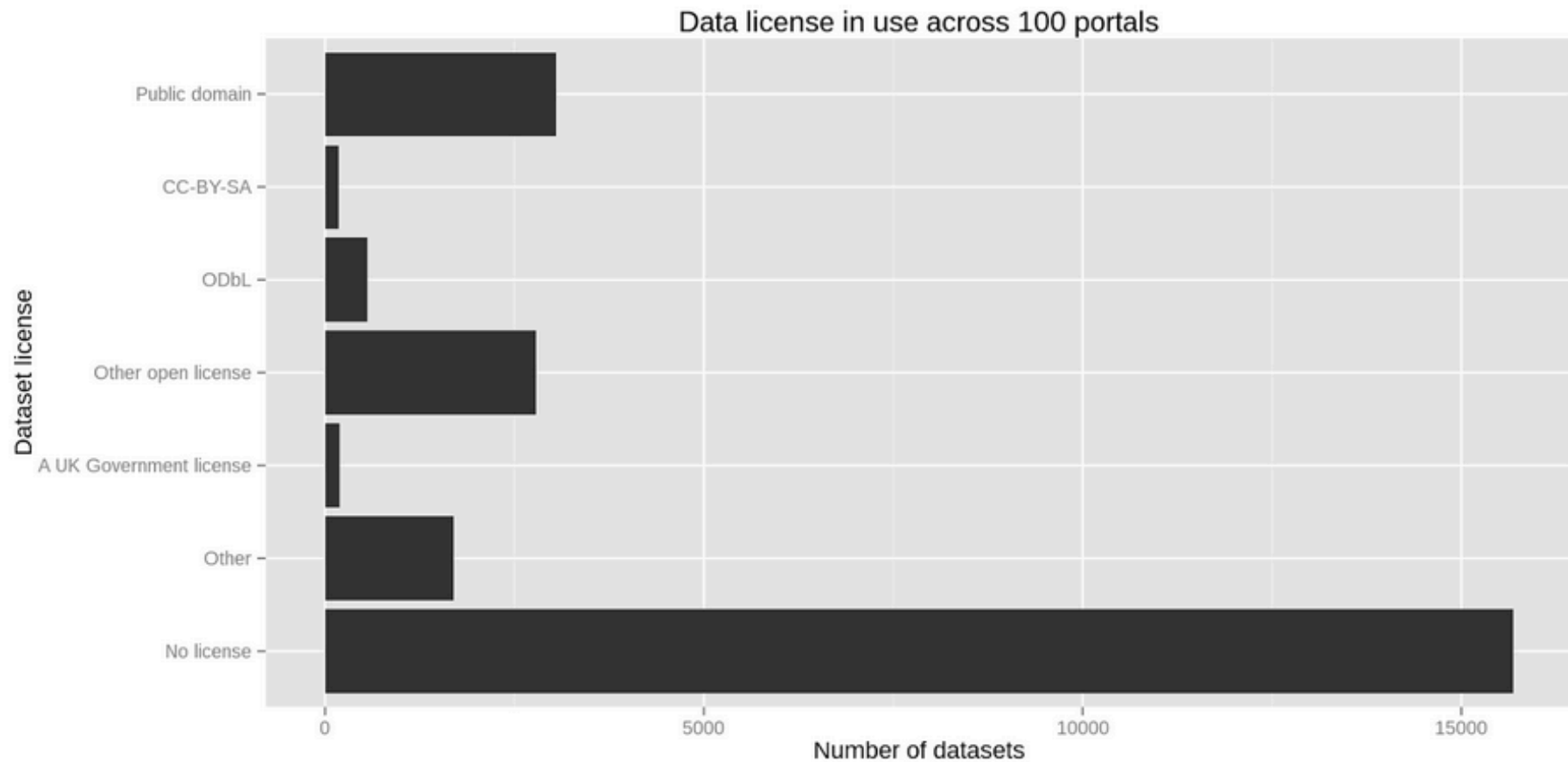
# Definition of Open (OKF)

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.

| | Open Definition (Open Knowledge Foundation) | OMB Memo, 2013 (The White House, Sylvia Burwell et al.) | Data.Gov.UK (Antonio Acuña) | "DBpedia: A Nucleus for a Web of Open Data" (Sören Auer et al.) | Open Data Institute (ODI) (Open Data Institute) | LinkedGov (LinkedGov) | McKinsey (James Manyika et al.) | Open Data Now (Joel Gurin) | Open Data Barometer (Tim Davies) | The World Bank (The World Bank) |
|---|---|---|---|---|---|---|---|---|---|---|
| Free | ✔ | ✔ | | ✔ | ✔ | | ✔ | | | |
| Negligible Cost | | | | | | | ✔ | | | |
| Publicly Available | ✔ | ✔ | | | ✔ | | ✔ | ✔ | | |
| Re-usable | ✔ | | ✔ | | ✔ | | | | | ✔ |
| Can be Redistributed | ✔ | | | ✔ | | | | | | ✔ |
| Non-exclusive (No Restrictions from copyright, patents, etc.) | ✔ | | | ✔ | ✔ | | | | ✔ | ✔ |
| Structured for Usability | | ✔ | ✔ | | | | ✔ | | ✔ | ✔ |
| Requires "Open" License | | | ✔ | | ✔ | ✔ | | | ✔ | ✔ |
| Non Personally Identifiable | | | | | | ✔ | | | | |
| Produced during business operation | | | | | | ✔ | | | | |
| Belongs to the Taxpayer (when not in violation of laws/privacy) | | | | | | ✔ | | | | |
| Accessible in Bulk | | | | | | | | | ✔ | |

GOVLAB

# Open data is hardly ever appropriately licensed.



Data license in use across 100 portals

Source: Thomas Levine
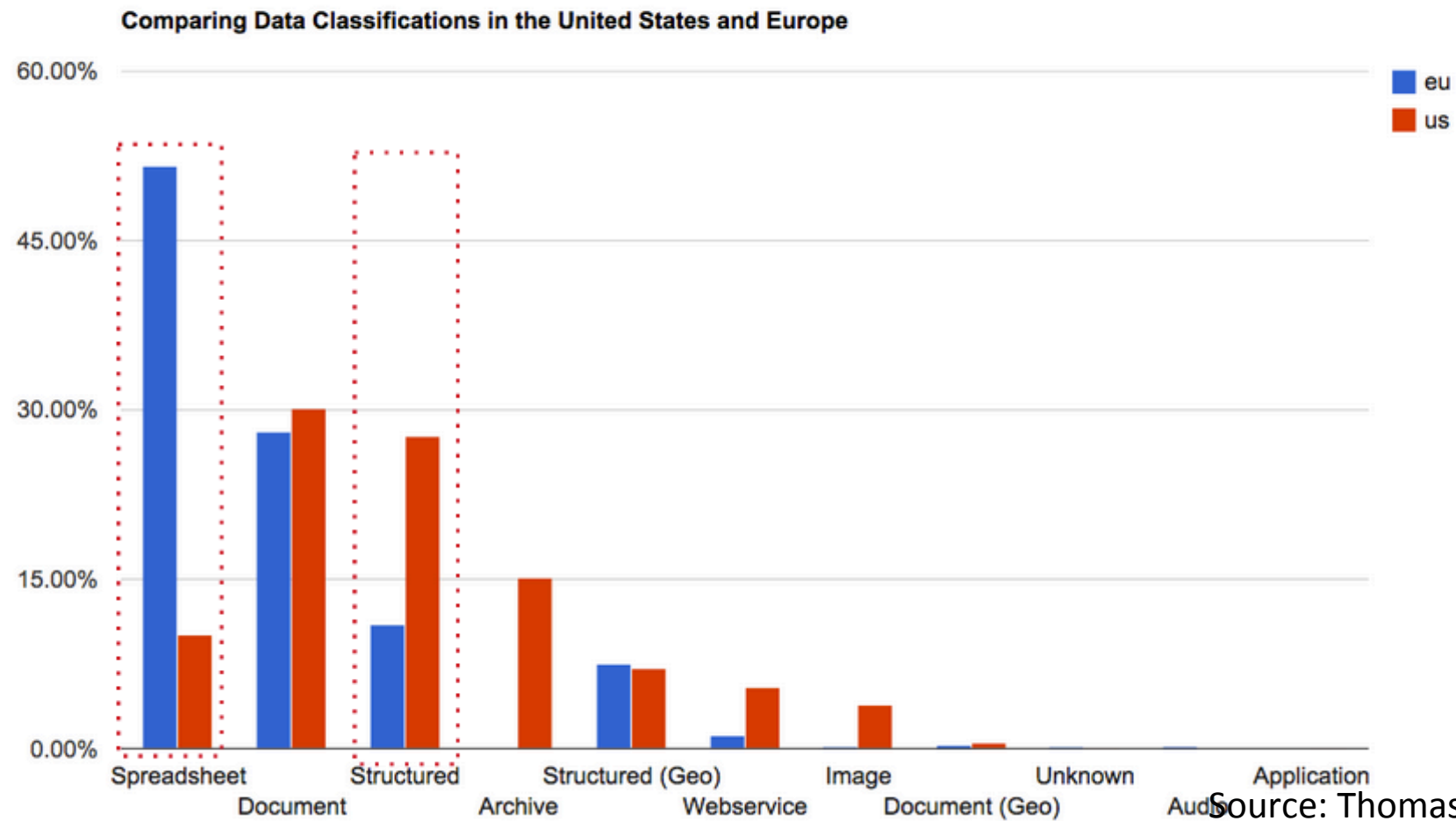
Hardly any datasets get updated.

Warning: **This result was actually based upon the wrong date field for updates. I personally think you'll find a similar result with the right date field being used.**

Source: Thomas Levine

# Open data is rarely structured.

**Comparing Data Classifications in the United States and Europe**



Legend:
- eu (blue)
- us (red)

Y-axis: 60.00%, 45.00%, 30.00%, 15.00%, 0.00%

X-axis categories: Spreadsheet, Document, Structured, Archive, Structured (Geo), Webservice, Image, Document (Geo), Unknown, Audio, Application

Source: Thomas Levine

# Today's mission

To move to phase 2 of publishing open data and solve some of the phase 1 problems.

# Publication phases

**Phase 1**: Get the data online, in some form. This will help with the trust and transparency and community building.

**Phase 2**: Increase the usability of the data by potentially publishing differently and keeping it up to date.

# Session 1
# The characteristics of data

# Exercise (part 1)

In your pre-training exercise, you were all asked to identify a dataset.

In your groups briefly discuss each others datasets and write down some key characteristics of each.

Also write the dataset title on a post-it, one per post-it.

# Types of Data

**Reference data**

**"things"**

**Transaction data**

**"stats involving things"**

# Exercise

Categorize your data into reference and transactional data.

If they are all in one category you have 2 minutes to add some new datasets to the empty category.

When done, put a "T" or and "R" on each dataset post-it.

Types of Data

Transaction data

"stats involving things"

Expenditure

Weather Consumption

Observation

Reference data

"things"

People  Facilities  Places
Books Buildings

# Update frequency

Static              In frequent updates              Frequent updates              Live

# Exercise

Categorize your data into **frequency of updates**

If they are all in one category you have 2 minutes to add some new datasets to the empty category/ies

Put a number on your post-its representing the frequency of updates.

  0 = static, 1= In frequent, 2 = Frequent, 3 = Live
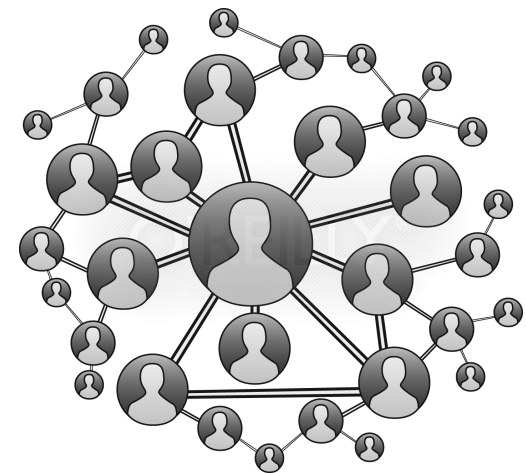
# Data Representations

## Tabular



## Hierarchical



## Network/Graph

# Exercise

Categorize your data into **tabular, hierarchical (tree) and graph (network)**

If they are all in one category you have 2 minutes to add some new datasets to the empty category.

Add the word "**tab**", "**tree**" or "**net**" to your post-its to represent the different structures.

# Justifications

Trust and Transparency | Enabling the economy

# One more

Categorize your data into **transparent** and **enabling**.

# Summing up

Do you have any obvious grouping of your datasets?

Is this reflective of the whole open data ecosystem?

# Exercise

Pick one "group" of datasets that share similar colours and come up with a data publication strategy for getting these datasets online and usable.

What are the publication requirements on the human publisher?

What are the requirements on potential users?

Thank-you