## ECONOMIC DEVELOPMENT

# DATA

### THE CASE FOR AN OPEN DATA PORTAL

## MARK WAINWRIGHT

ADB, MANILA, MAY 2013

Open Knowledge Foundation

The Open Knowledge Foundation aims to open up knowledge around the world to make it accessible and useful. It is home to many activities, including an Open Government working group and a consulting division. OKF is a non-profit organisation registered in the UK.
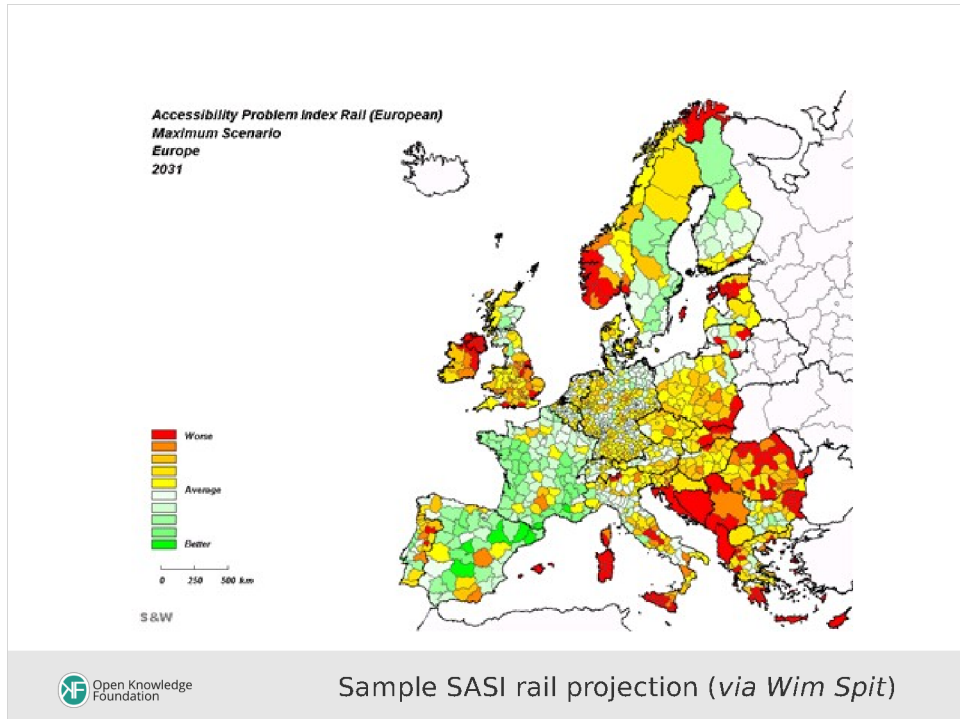
**Passenger and freight turnover data, all modes of transport**
*Source: Vakulchuk & Irnazarov: CAREC case study*

| NSO | Period | Completeness | Region (oblast) level | Gaps |
|---|---|---|---|---|
| Kazakhstan | 2003 | High | Yes | Cross-border transactions |
| Uzbekistan | 2009 | Medium | Yes | Cross-border transactions Incomplete |
| Kyrgyzstan | 1990 | Medium | Yes | Cross-border transactions Incomplete |
| Tajikistan | 2000 | Low | No | Micro-level data |

Open Knowledge Foundation

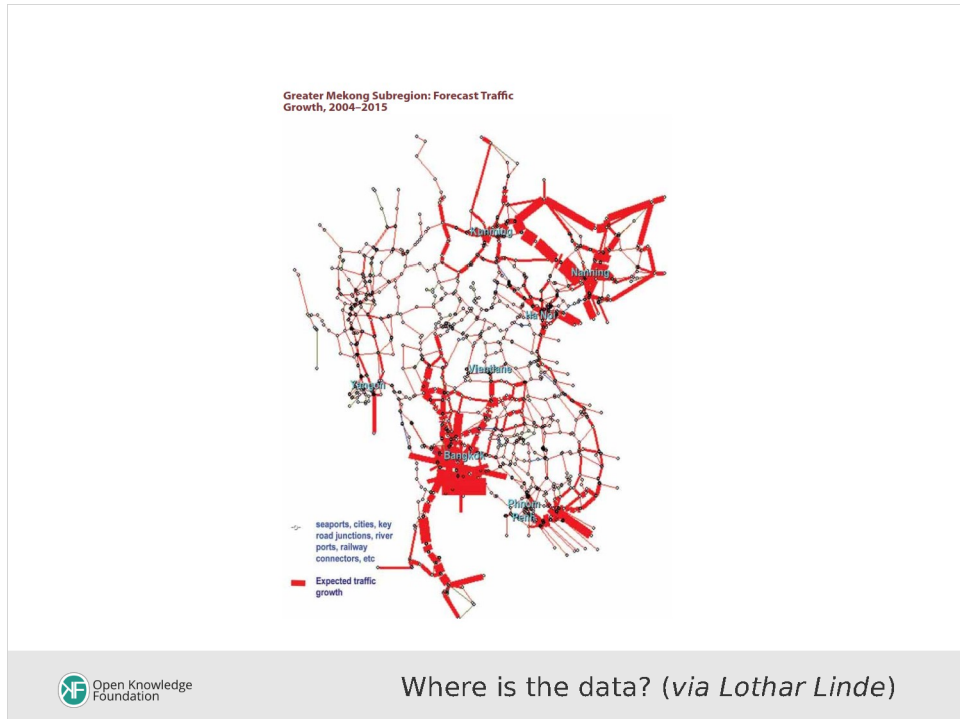Missing transport data in the Stans

First, the bad news - a lot of data that one would like is missing, incomplete, or inconsistent, or is collected in incommensurate ways across national borders.

Accessibility Problem Index Rail (European)
Maximum Scenario
Europe
2031

Worse

Average

Better

0    250    500 km

S&W

Open Knowledge Foundation
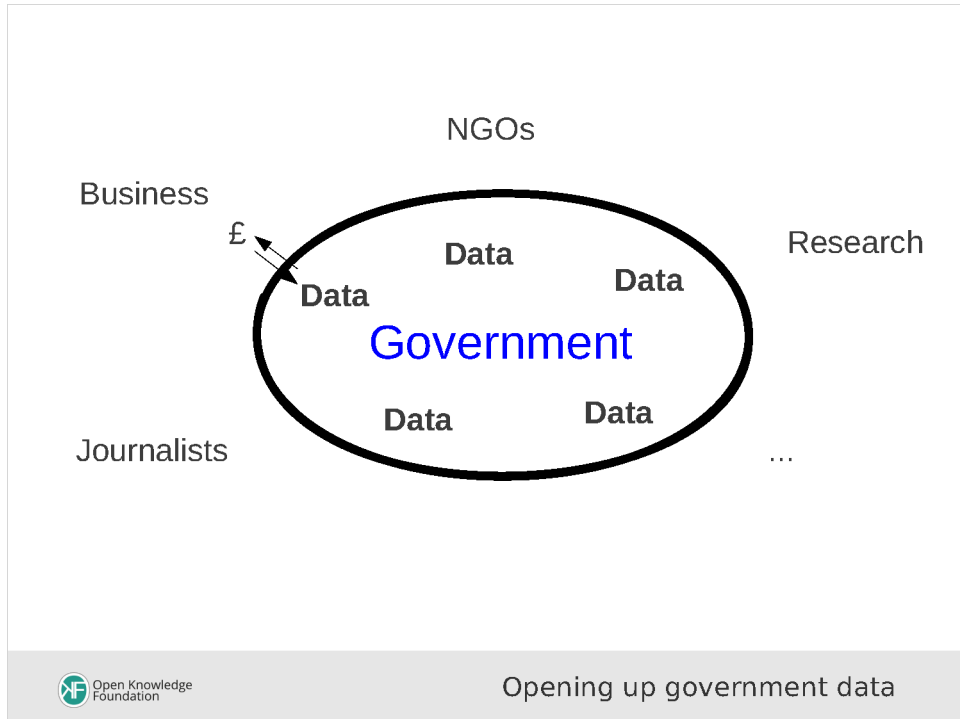
Sample SASI rail projection (*via Wim Spit*)

This is a projection (of rail accessibility) in Europe, from the SASI model, which models socio-economic and spatial impact of different transport policy options. It's a rich model but uses base data of many variables (regional GDP, population by age and sex, labour force participation, migration, etc) over 1330 districts, as well as transport network data. This kind of analysis can provide valuable insights and guidance but requires a consistent data framework across the region.

K Prasad: Agent-Based Model for the Evaluation of Aid for Trade Infrastructure Investments

Open Knowledge Foundation

South Asia: districts and tiles

While many datasets are collected at district level, some international datasets are divided into approximately square tiles covering 1 degree of longitude and latitude, particularly where the data is compiled or estimated from satellite images. This provides a further challenge when data from both types of source needs to be integrated. It's essential for data to be geo-coded so that data from different sources can be overlaid. SASEC used tile data, together with data on transport links between tiles, to model local effects of improving transport links to reduce market friction in South Asia, shown here with both tiles and districts.

Greater Mekong Subregion: Forecast Traffic Growth, 2004–2015

seaports, cities, key road junctions, river ports, railway connectors, etc

Expected traffic growth
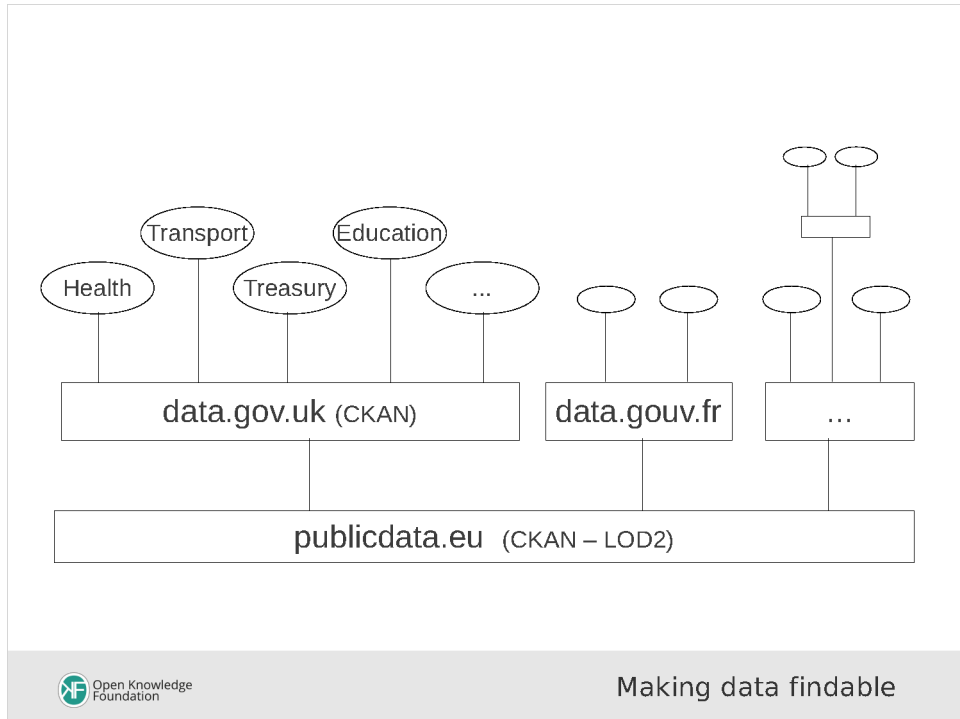
Where is the data? (*via Lothar Linde*)

Speaking of transport data, this graphic was shown me by Lothar Linde of the GMS project. It shows projected traffic flows in the GMS region (thicker lines represent more traffic; Bangkok is seen to be a hub). The graphic appeared in an ADB report. Lothar wanted to use the underlying data, but he was unable to trace it despite several months of searching. Solving the problem of data that goes missing or can't be found is both far more urgent than gathering new data, and far easier.

NGOs

Business

£

**Data**

**Data**

**Data**

Research

Government

**Data**

**Data**

Journalists

...

Open Knowledge Foundation

Opening up government data

Governments necessarily collect a lot of data, which traditionally was kept secret (or perhaps monetised) but is increasingly being released openly and freely, with Open Government Data (OGD) initiatives across the world, though very few in Asia. Benefits of OGD include transparency, increased participation and engagement, and socio-economic value. Data can be used to power new types of services and generate value and economic activity. There are also obvious benefits to research. Increasingly many governments are being convinced that the benefits far outweigh the cost of lost revenue.

Open Knowledge Foundation

Thomas Herndon is a graduate student who, in a famous recent story, found a crucial error in the data behind an influential paper of Reinhart and Rogoff (2010). To get hold of the spreadsheet he only needed to e-mail and ask for it, seemingly a small barrier. But if their data had been provided for download on the internet when the paper was published, would it have taken 3 years for the error to come to light?

Making data findable

Even if data is available for download and can be freely re-used, it's not much use unless people can find it. UK government data is released by hundreds of different departments and bodies, each with their own website - how to find it? The answer is in data.gov.uk, a data catalogue or 'portal' where all of it is registered, and can be found easily with searching, filtering, etc. data.gov.uk currently knows of nearly 10,000 datasets. publicdata.eu is an experimental EU-wide catalogue that 'harvests' data records from a number of different national and local portals. Though these portals are principally used as a catalogue, a data portal can also be used to store the data itself. We'll look in more depth at later at CKAN, an example data portal.

data.gov.uk was mentioned above. Note the number of datasets, the 'data' tab (giving more advanced search options), and the 'Request data' tab, allowing data users to ask for particular datasets to be made open if they are not already on the site.

Greater Mekong Subregion: Forecast Traffic Growth, 2004–2015

If ADB had had a data portal and a policy of publishing data used in reports etc, the data behind this visualisation would probably not have gone missing. The example shows that it is valuable to publish model data as well as baseline data.

The EU is another major publisher of data. Note that this portal is for the EU's own data - it is not the same as publicdata.eu (mentioned above) which aggregates data records from member states. A number of EU offices and bodies publish data, including their own administrative data, but at present the great majority of the data that can be found here is from Eurostat, the EU statistics agency.

Among the data at data.gov.uk is departmental spending data - every department is supposed to publish this. The Cabinet Office team in charge of data.gov.uk put together this dashboard to show how well departments were doing. This sign of interest acted as a nudge to a number of departments that had not previously published any spending data. An Open Data policy is most effective when, as here, it comes from a high level - as with the recent US open data directive, issued by the President.

Source: Roland-Holst & Sugiyarto

Understanding models

An example of what a 'dashboard' interface to model data might look like, presenting it in a useful form to policy makers. (Of course, the data needs to be available first!) This example uses commercial data visualisation tools.

The default theme for CKAN 2.0, the free, open-source data portal software maintained by the Open Knowledge Foundation. CKAN has been developed since 2007 and now powers portals such as the UK and EU ones described above. By chance, the release of CKAN 2.0 was announced last week. The following slides show a very small number of its features, in particular in the area of geodata, whose importance was stressed above.

Data can be added to CKAN in a number of ways, including harvesting or ingesting in bulk from existing information systems. This shows part of the user interface for adding a dataset by hand. Various metadata is added such as title, description, etc (the metadata schema can be changed). Tags help users find the dataset when searching. Data should have enough metadata that users can find the data they need and know what data they are dealing with.

data.gov, the US government data portal, is due to be relaunched (using CKAN) in a few weeks. This prototype has over 50,000 datasets. As well as searching by keywords, tags, publishing department, etc, note the map which allows searching by geographic extent where available.

Drawing a rectangle on the map applies a filter to show only datasets that intersect with it. Note the much reduced number of results. Further searching can be carried out within the results.

The description of the first dataset on the previous slide. Note the dataset extent map.

Along with information about data found, CKAN can show previews of various kinds of data. This shows the ability to preview WMS geodata, including selecting different layers to overlaly.

## Key Action Points, 6 to 12 months

**6 months**

- Develop and approve study proposal
- Interactions at high policy level for opportunity scenarios
- Scope out an ADB publishing platform to make existing data discoverable for pilot regions
- Policy meeting for endorsement

**12 months**

- Pick pilot region with high-level policy backing
- Get funding and resources in place
- Pilot government data portal in countries within pilot region
- Design a user friendly knowledge tool
- Decide on an institutional framework for ADB data publishing platform

A reprise of the key action points seen in the previous presentation.