# Toward a Frictionless Data Future

PRESENTED BY

## Dan Fowler

**Developer Advocate**

daniel.fowler@okfn.org (@danfowler & @okfnlabs)

OPEN KNOWLEDGE INTERNATIONAL
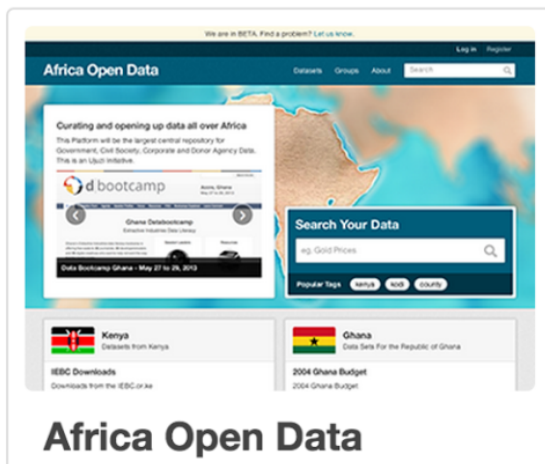
FRICTIONLESS DATA
STANDARDS AND TOOLING

# Who we are

International non-profit founded in 2004

- **Vision**
  - A world where **open knowledge is ubiquitous**
- **Mission**
  - **Open up** all essential public-interest information
  - **Build communities, tools, and skills** to empower individuals to use open information to drive change

OPEN KNOWLEDGE INTERNATIONAL

OPEN KNOWLEDGE

# CKAN -> Frictionless Data

- Our work on Frictionless Data is borne out of what we've learnt from building and deploying CKAN
- Working with publishers and data publication flows all around the world has highlighted the "friction" involved in working with data



Africa Open Data



Buenos Aires Data



City of Ottawa Open

# "Tidy" Data

## Framework for prepping data for analysis

*It is often said that 80% of data analysis is spent on the process of cleaning and preparing the data (Dasu and Johnson 2003). Data preparation is not just a first step, but must be repeated many over the course of analysis as new problems come to light or new data is collected. [...] Part of the challenge is the breadth of activities it encompasses: from outlier checking, to date parsing, to missing value imputation.*

*Hadley Wickham, Tidy Data (2014)*
*https://www.jstatsoft.org/article/view/v059i10*

FRICTIONLESS DATA
STANDARDS AND TOOLING

## Frictionless Data:

- Lightweight extensible **specifications** for "packaging" datasets

    Current focus: **Tabular Data**

- **Integrations** for loading datasets into tools and platforms relevant to researchers
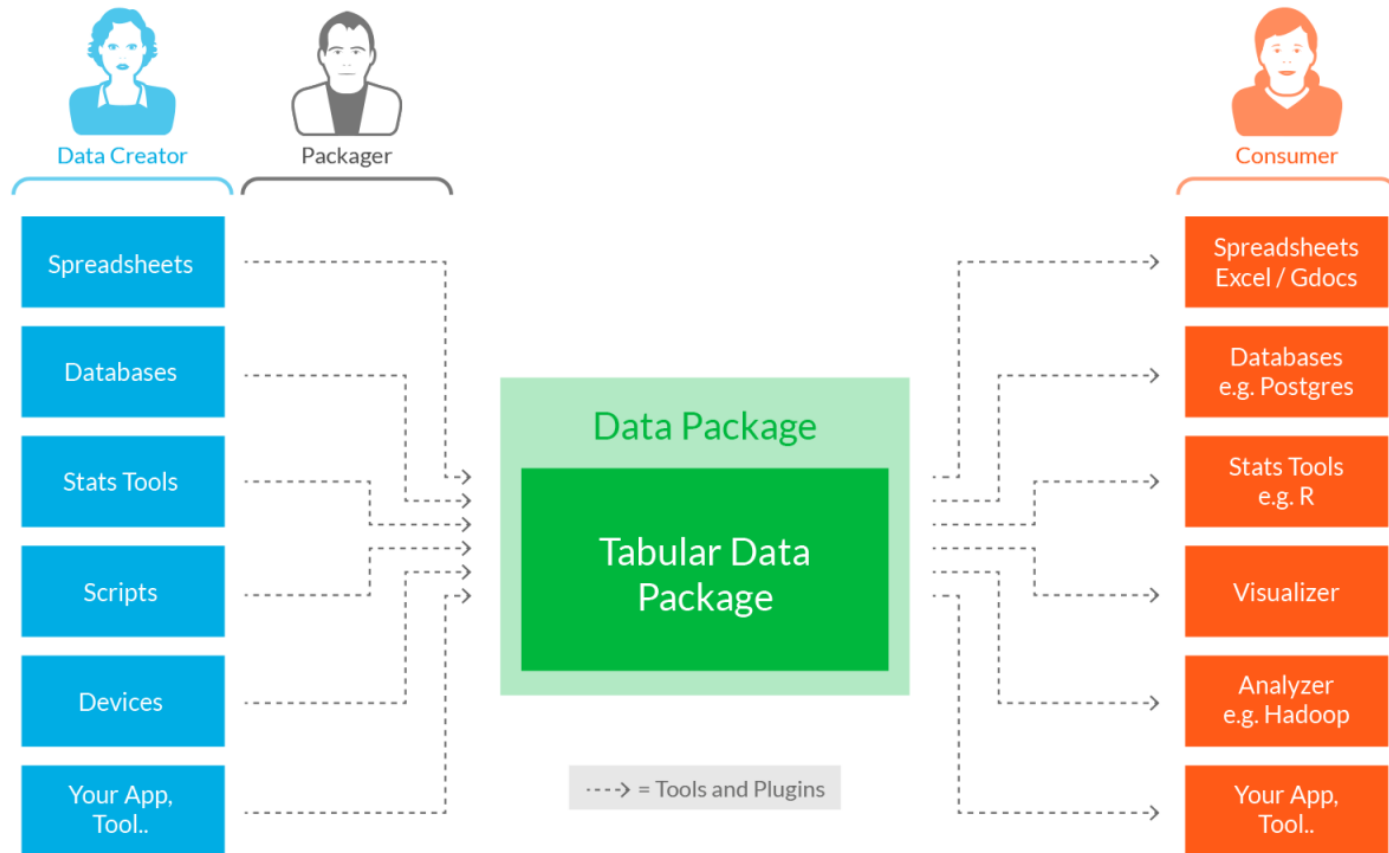
## Aims:

- Introduce a significant, *measurable* improvement in how research data is shared, consumed, and analyzed.
- Make it easier to maintain and improve data **quality**.

OPEN KNOWLEDGE

# Integrations



Data Creator

Packager

Consumer

Spreadsheets

Databases

Stats Tools

Scripts

Devices

Your App, Tool..

Data Package

Tabular Data Package

- - -> = Tools and Plugins

Spreadsheets
Excel / Gdocs

Databases
e.g. Postgres

Stats Tools
e.g. R

Visualizer

Analyzer
e.g. Hadoop

Your App,
Tool..

OPEN KNOWLEDGE

# Data Containerization

# Data

## data.csv

```
Date,      Age,  Primary Treatment Outcome
1/1/2012,  29,   PARTIAL
4/5/2012,  32,   COMPLETE
6/2/2012,  18,   PARTIAL
8/12/2012, 20,   PARTIAL
19/1/2012, 25,   COMPLETE
```

OPEN KNOWLEDGE

# Data (Packaged)

datapackage.json

```json
{
  "name": "my-data",
  "title": "My Data",
  "sources": ["http://example.com/data.csv"],
  "license" : "ODC-PDDL-1.0",
  "resources": [{
    "path": "data.csv",
    "format": "csv",
    "dialect": {
      "delimiter": ",",
      "header": true
    },
    "schema": {
      ...
    }
  }]
}
```

OPEN KNOWLEDGE

# JSON Table Schema

```json
"schema": {
  "fields": [
    {
      "name":"Date", "type":"date", "format":"fmt:%d/%m/%Y",
      "constraints": {"required": "true"}
    },
    {
      "name":"Age", "type":"integer",
      "missingValue": "NULL"
    },
    {
      "name":"Primary Therapy Outcome",
      "description": "Whether primary treatment success",
      "type": "string",
      "constraints": {"enum": ["PARTIAL","COMPLETE"]}
    }
  ]
}
```

# Tabular Data Package

| Tabular Data Package | | |
|---|---|---|
| Data Package | JSON Table Schema | CSV Dialect Description |

http://specs.frictionlessdata.io/csv-dialect/

http://specs.frictionlessdata.io/data-packages/

http://specs.frictionlessdata.io/json-table-schema/

http://specs.frictionlessdata.io/tabular-data-package/

# **Benefits of Approach**

# Focus on Extensibility

# Tools

# Solid Foundation for Tools



**Tool**
e.g. validation via "GoodTables"

**Tool**
e.g. direct import into R, Pandas, SQL, etc.

**Tabular Data Package**

| Data Package | JSON Table Schema | CSV |

OPEN KNOWLEDGE

# Good Tables

- GoodTables is a tool built by Open Knowledge International for data source validation.
- Once you have defined a schema using JSON Table Schema, you can make sure your CSV file is valid against the schema

# Platform Integrations

# Partners

# Partners

Western Pennsylvania
Regional Data Center

DM4T
EPSRC Data Management for TEDDINET using Semantic Technologies

ROpenSci

Dataship

Pacific Northwest
NATIONAL LABORATORY

UNIVERSITY OF CAMBRIDGE

# Seeking Pilot and Technical Partnerships

- Project overview: http://frictionlessdata.io

- Partners: http://frictionlessdata.io/partners/

- Case Studies: http://frictionlessdata.io/case-studies/

- Specifications: http://specs.frictionlessdata.io/

- **Newsletter**: http://frictionlessdata.io/get-involved/#newsletter

# Questions?