

8th International
Conference on
BIG DATA
& Data Science for Official Statistics

BILBAO 2024

Informing Climate Change and
Sustainable Development Policies
with Integrated Data

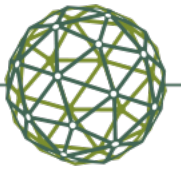
BILBAO, SPAIN | **10-14 JUNE 2024** | **#UNBigData2024**

UN Data and the Use of Data Commons Infrastructure



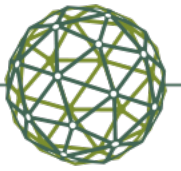
UN Data

Powered by Google's Data Commons



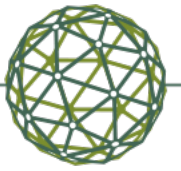
The UN Data modernization project

- Response to UN Secretary-General's Data Strategy
- Informed by the Roadmap for Innovating UN Data and Statistics
- Endorsed by UN Secretary-General's Executive Committee



Overall objectives

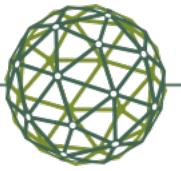
- Increase visibility and access to authoritative sources of statistical data and metadata
- Improve search and analytic capabilities for policy and decision makers
- Enable interoperability of statistical data from across the UN system, Member States and other partner organizations
- Enhance data value through meaningful interlinkages across global, regional and national data portals



Concrete application of CDIF* principles

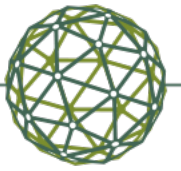
- Harness the UN's authority, credibility, and name recognition to unify efforts and collaborate across all stages of the data life cycle
- Ensure trust in data through data quality and adherence to standards
- Develop capacity of all stakeholders to both contribute and use a common, modern web infrastructure

* Cross-Domain Interoperability Framework <https://doi.org/10.5281/zenodo.11236871>



Building on existing standards, infrastructure and communities of practice

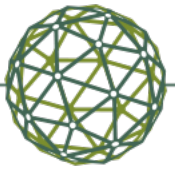
- Collaborate with domain experts in data modeling and integration tasks to validate transformations and mappings
- Use .Stat as dimensional data repository, leveraging SDMX standards
- Deploy on UN Global Platform infrastructure
- Engage with statistical community to collaboratively manage data integration process



Towards a distributed UNdata Knowledge Graph

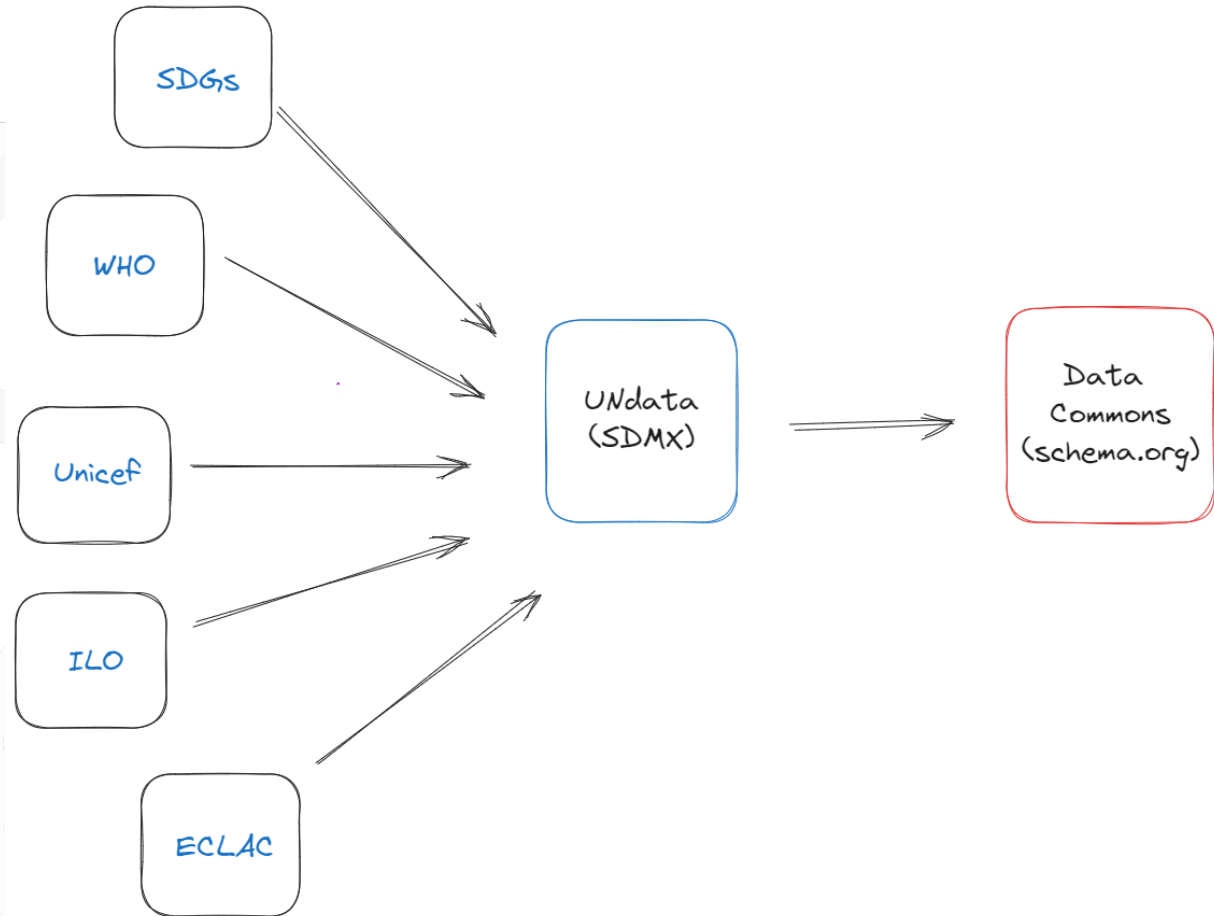
The objective is to capture the concepts and relationships required to:

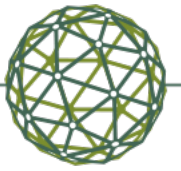
- Establish explicit and implicit links to external resources, thus making data more easily findable, searchable, and usable.
- Build applications that efficiently access related data across multiple domains using linked open data techniques
- Generate insights by reasoning over complex relationships.
- Incrementally add new data and evolve the data schema to accommodate new data types and new use cases.



Objective

Integrate datasets from multiple, heterogeneous sources, by converting them to a common UN Data schema, and ingest them into the Data Commons Knowledge Graph (based on schema.org)

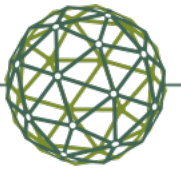




Towards a distributed UNdata Knowledge Graph

The objective is to capture the concepts and relationships required to:

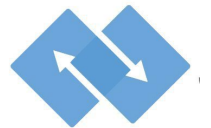
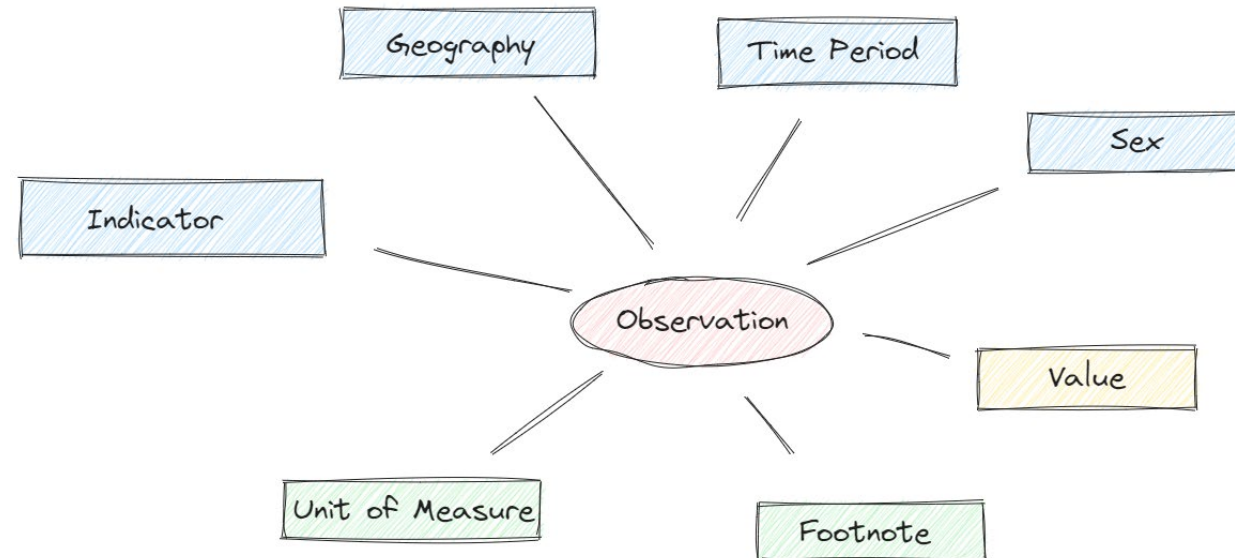
- Establish explicit and implicit links to external resources, thus making data more easily findable, searchable, and usable.
- Build applications that efficiently access related data across multiple domains using linked open data techniques
- Generate insights by reasoning over complex relationships.
- Incrementally add new data and evolve the data schema to accommodate new data types and new use cases.



Data integration problem

- Equivalent concepts (entities) are assigned different identifiers in different databases or vocabularies
- Different agencies use different names for the same real-world entities
- A significant part of the work consists in aligning terms that refer to the same concept across datasets.
- Mapping rules are often hidden or not documented at all
 - Describe the correspondences between classification schemes
 - Tracking how classification items have been created, split, merged, or removed from active use

Concept schemas



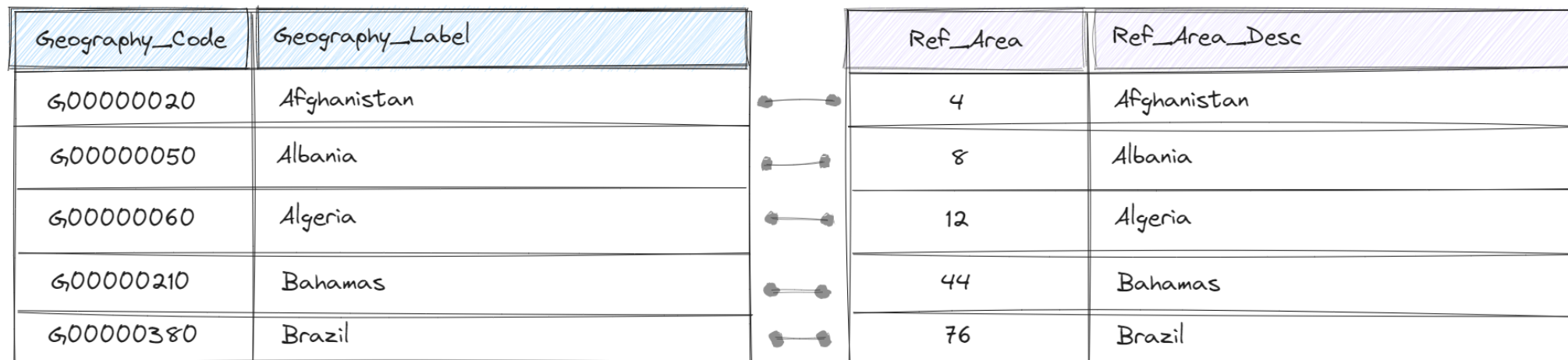
Enumerations

Sex_Code	Sex_Label
_T	Total
F	Female
M	Male

Geography_Code	Geography_Label
G00000020	Afghanistan
G00000050	Albania
G00000060	Algeria
G00000210	Bahamas
G00000380	Brazil



Mappings between source entities and UN Data entities



Data pipelines



Source

Source converted to "Proto SDMX"

Source dataset

Indicator	Geography	Time Period	Sex	Value	Unit of Measure	Footnote
Life expectancy	Bhutan	2023	Total	74.00	Percent	Large group
Life expectancy	Bhutan	2023	Male	71.00	Percent	Large group
Life expectancy	Bhutan	2023	Female	77.00	Percent	Large group
Life expectancy	Bhutan	2022	Total	73.00	Percent	Large group
Life expectancy	Bhutan	2022	Male	70.00	Percent	Large group
Life expectancy	Bhutan	2022	Female	76.00	Percent	Large group
Life expectancy	Vietnam	2023	Female	75.00	Percent	Large group
Life expectancy	Vietnam	2023	Total	74.00	Percent	Large group
Life expectancy	Vietnam	2023	Male	73.00	Percent	Large group
Life expectancy	Vietnam	2022	Female	74.00	Percent	Large group
Life expectancy	Vietnam	2022	Total	73.00	Percent	Large group
Life expectancy	Vietnam	2022	Male	72.00	Percent	Large group

Source enumerations

Geography Code	Geography Label	Sex Code	Sex Label
v	Aggregation	SEX_T	Total
f	Albania	SEX_M	Male
24	Angola	SEX_F	Female
44	Bahrain		
74	Brazil		
83	Bolivia		

Source schema

```

graph TD
    Observation((Observation)) --- Indicator[Indicator]
    Observation --- Geography[Geography]
    Observation --- TimePeriod[Time Period]
    Observation --- Sex[Sex]
    Observation --- Value[Value]
    Observation --- UnitOfMeasure[Unit of Measure]
    Observation --- Footnote[Footnote]
  
```



curation

Mappings

UN Data

Geography Code	Geography Code
4000000020	v
4000000030	f
4000000040	24
4000000050	44
4000000060	74
4000000070	83

Sex Code

Sex Code	Sex Code
_T	SEX_T
M	SEX_M
F	SEX_F



curation

Mappings

UN Data

Geography Code	Geography Code
4000000020	v
4000000030	f
4000000040	24
4000000050	44
4000000060	74
4000000070	83

Sex Code

Sex Code	Sex Code
_T	SEX_T
M	SEX_M
F	SEX_F

Harmonized dataset

UN Data dataset

Indicator	Geography	Time Period	Sex	Value	Unit of Measure	Footnote
Life expectancy	Bhutan	2023	Total	74.00	Percent	Large group
Life expectancy	Bhutan	2023	Male	71.00	Percent	Large group
Life expectancy	Bhutan	2023	Female	77.00	Percent	Large group
Life expectancy	Bhutan	2022	Total	73.00	Percent	Large group
Life expectancy	Bhutan	2022	Male	70.00	Percent	Large group
Life expectancy	Bhutan	2022	Female	76.00	Percent	Large group
Life expectancy	Vietnam	2023	Female	75.00	Percent	Large group
Life expectancy	Vietnam	2023	Total	74.00	Percent	Large group
Life expectancy	Vietnam	2023	Male	73.00	Percent	Large group
Life expectancy	Vietnam	2022	Female	74.00	Percent	Large group
Life expectancy	Vietnam	2022	Total	73.00	Percent	Large group
Life expectancy	Vietnam	2022	Male	72.00	Percent	Large group

UN Data enumerations

Geography Code	Geography Label	Sex Code	Sex Label
v	Aggregation	SEX_T	Total
f	Albania	SEX_M	Male
24	Angola	SEX_F	Female
44	Bahrain		
74	Brazil		
83	Bolivia		

UN Data schema

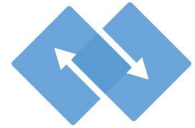
```

graph TD
    Observation((Observation)) --- Indicator[Indicator]
    Observation --- Geography[Geography]
    Observation --- TimePeriod[Time Period]
    Observation --- Sex[Sex]
    Observation --- Value[Value]
    Observation --- UnitOfMeasure[Unit of Measure]
    Observation --- Footnote[Footnote]
  
```



MCF FILE

Contrast between SDMX Information Model and Data Commons schema

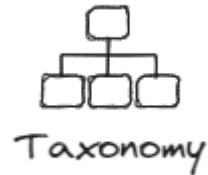


- Specifically designed for statistical data, with a more tabular, dimensional focus.
- More rigid but specific, focusing on dimensions, attributes, and measures.
- Uses StructureMaps and ComponentMaps to describe how data should be transformed or related.
- Better suited for statistical data where dimensions and measures are predefined.

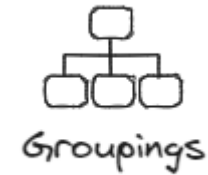


- Based on the Schema.org model, with roots in knowledge representation systems.
- Aims for a more flexible, verbose base layer, allowing various kinds of relationships and attributes
- Provides different APIs (Node, SPARQL, DCGET) for various views, including a time-series view.
- Built as a knowledge graph, making it more suitable for capturing a wide range of relationships among diverse entities

Additional data structures



- Taxonomy of UN Data themes and sub-themes

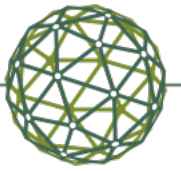


- Definition of statistical variable groupings

Enriched metadata

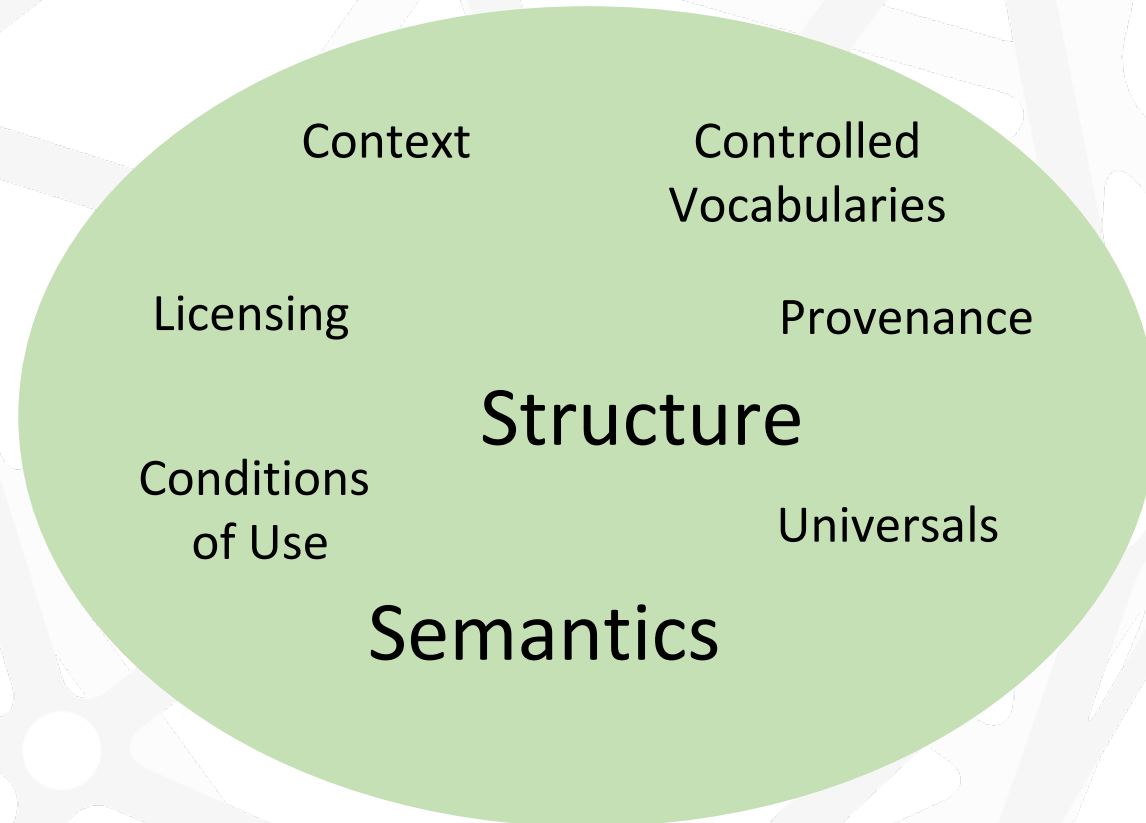
Additional metadata is captured for each indicator or sv-grouping:

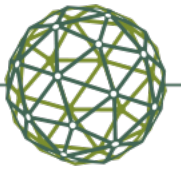
- Metadata link
- Revised sv-group description



Functions of Cross-Domain FAIR implementation

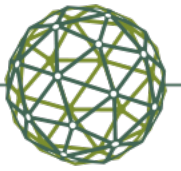
- Discovery
- Assessment
- Access
- Integration
- Packaging





Alignment with CDIF recommendations

- Re-express domain metadata (e.g., SDMX) into other common, integration-ready data descriptions (e.g., schema.org).
- Manage data at the finest granularity possible to facilitate composability
- Use persistent, de-referenceable identifiers at all levels of granularity



Making UN Data AI-ready

- Data profiling for **Discovery, Assessment** and **Access**
- Mappings of concepts and code lists for **Integration**
- Supplementary metadata for **Context** and **Provenance** traceability (e.g., edited indicator descriptions, mapping to UN thematic areas, definition of variable groupings, methodological notes...)



Department of Economic and Social Affairs
Statistics



Home

Countries / Areas

Thematic Areas

SDGs

Data Partners

Search



UN Data

Powered by Google's Data Commons

How much food goes wasted around the world?



Introducing UN Data, a powerful tool for extracting insights from data available across the UN system. Search and explore high-quality datasets and digital public goods right at your fingertips to empower and expedite evidence-based decision-making.

<https://unstats.un.org/UNSDWebsite/undatacommons/>



UN **data**

Powered by Google's Data Commons



**United
Nations**

Department of
Economic and
Social Affairs