



**United
Nations**

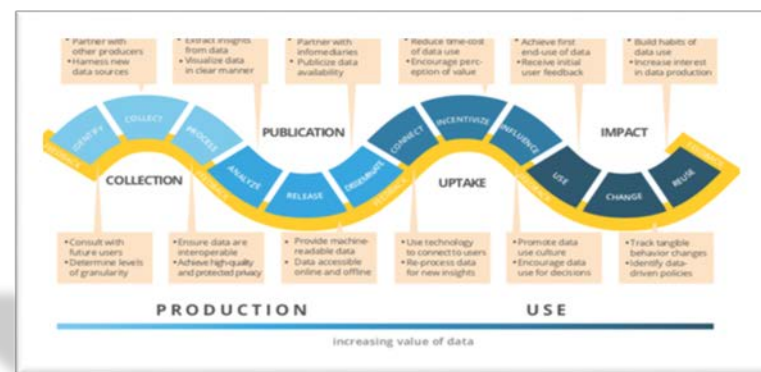
Department of Economic and Social Affairs
Statistics

Introduction to SDMX data modeling

Interoperability

- Ability to seamlessly share, join, cross-analyse, exchange and re-use data produced from different sources, and at different times, to provide richer information for improved decision making
- It is a crucial characteristic of good quality data and of effective data management systems

Interoperability should be understood along the whole “data value chain”, from **collection** to **use**



Interoperability and data modelling

- Interoperability is highly dependent on data and metadata modelling decisions and practices
- The same information content is often represented in variety of ways across different systems and organizations.
- There is usually no single “right” way of representing information
 - Some data structures are better suited for managing transactional processes (e.g., capturing data from a survey or maintaining a civil registration database)
 - Others are better suited for analyzing and communicating data to users (e.g., for the creation of data visualizations in a monitoring dashboard).

What is data modeling?

- A process focused on:
 1. Clearly and unambiguously **identifying things** that a dataset aims to capture
 2. Selecting the key properties that should be captured to **describe those things** in a meaningful way
 3. Deciding **how things relate** to each other
 4. Deciding how this information should be **formally codified**
- *This is the essence of the Entity-Relationship model, which underlies most modern database management systems and applications*

Examples

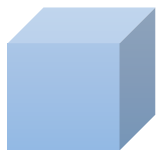
- The content of a dataset may refer to **entities** such as “city”, “person”, or “activity”,
- These entities may be usefully described with **attributes** like “name”, “age”, or “industry”.
- In a specific application, it could be useful to capture **relationships** among entities and attributes, e.g., the fact that
 - one or more persons may live in a city,
 - a person may be employed in one or more types of activity ...

Canonical data and metadata models

- Models that follow specific standardized patterns, making them highly **reusable** and conducive to data sharing.
- Provide a **common template** to which different datasets can be mapped
- Help develop a common understanding of how the various components of a dataset relate to each other and to the components of other datasets
- Reduce the number of transformations that user applications need to perform on their own to integrate the data from those sources

Standardization is not for free

- The underlying principle is **to hide from user the internal complexity of the operational data models** (e.g., which are optimized to avoid data redundancy and ensure data consistency validations), so they can concentrate on using data rather than spending time trying to understand the intricacies of internal data structures
- Data **providers need to take responsibility for mapping the data** from its original, operational structures, into commonly agreed presentations for dissemination and distribution purposes
- This may entail the need to undertake so-called “Extract-Transform-Load”, or ETL, procedures, **hidden from the view of users**



The multi-dimensional ‘data cube’ model

- Presents all relevant data about a population of interest in a **simple, self-contained** tabular view
- Each data point is characterized by
 - **Measures**: Observed values on one or more variables interest
 - **Dimensions**: A set of uniquely identifying characteristics
 - **Attributes**: A set of additional characteristics that further describe it

Domains of dimensions, attributes and measures

- Each dimension, measure and attribute encapsulates **a concept**
- Concepts can be:
 - drawn from a code list (for e.g., "country ISO code")
 - required to adhere to a specific data format (e.g., "YYYY" for years)
 - required to be contained within a specific range of values (e.g., "numerical values between 0 and 1").
 - drawn from a type of values (e.g., "text")



**United
Nations**

Department of Economic and Social Affairs
Statistics

The SDMX information model

Figures vs data

Number of Tourist establishments Italy, Annual data			
Indicator Time	A100 Hotels and similar	B010 Tourist Campsites	B020 Holiday dwellings
2002A00	33411	2374	61479
2003A00	33480	2530	58526
2004A00	33518	2529	56586
2005A00	33527	2411	68385
2006A00	33768	2510	68376
2007A00	34058	2587	61810

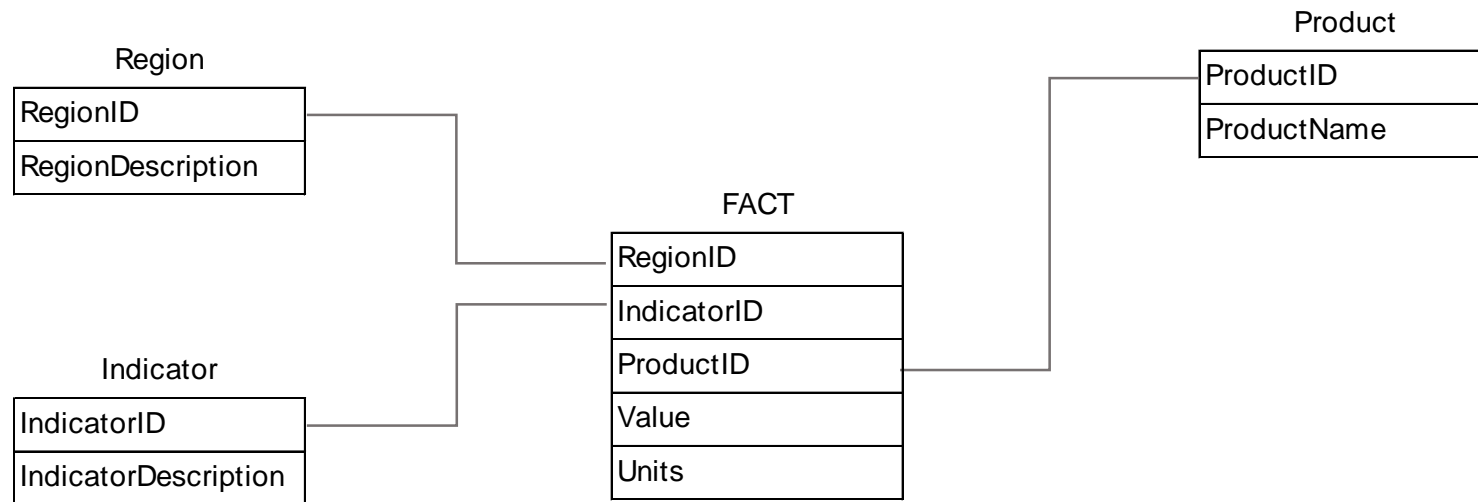
- Figures by themselves are meaningless.
- For data to be usable, it must be properly described. The descriptions let users know what the data actually represents.

Developing a Data Model for Data Exchange

- Data model is developed to provide descriptions for all relevant characteristics of the data to be exchanged
- In some aspects similar to developing a relational database
- In SDMX, data model is represented by a Data Structure Definition (DSD).
- The “shape” of SDMX DSD is roughly similar to star schema.
- To design a DSD, we first need to find *concepts* that identify and describe our data.

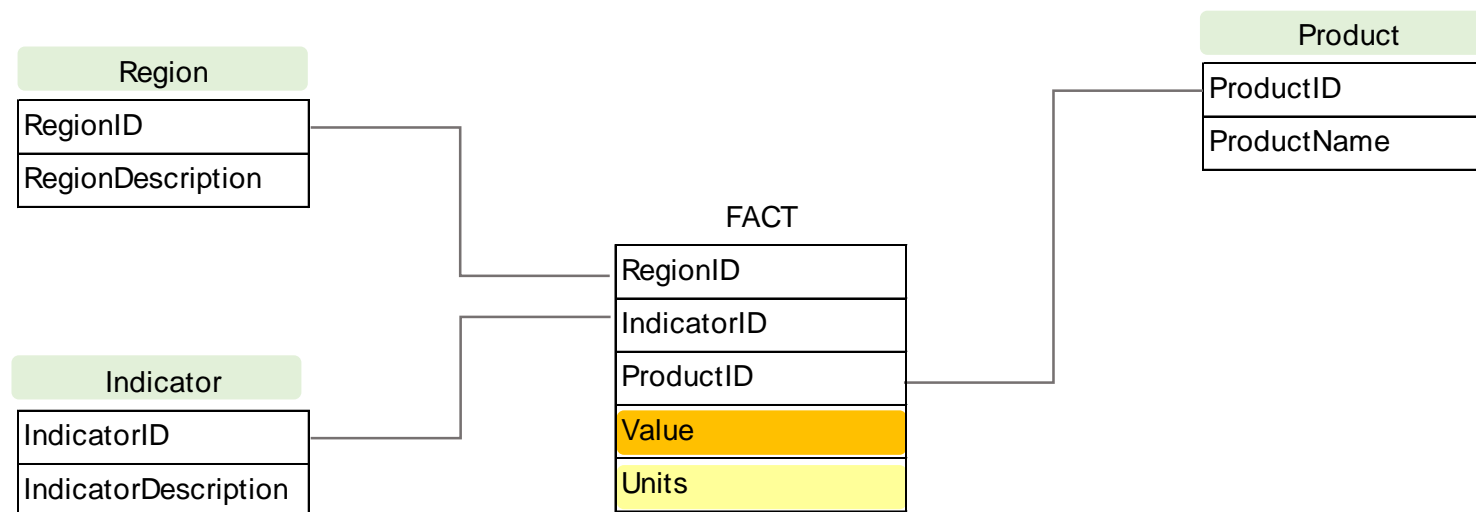
Developing a Data Model for Dissemination and Exchange

- This task is similar to developing a relational database
- In SDMX, a data model is represented by a **Data Structure Definition (DSD)**
- The “shape” of SDMX DSD is roughly similar to a “**star schema**”



Concept

- To design a DSD, we first need to find **concepts** that describe all relevant characteristics of our data



Identifying Concepts

Ref. Area

Indicator

Time Period

Unit Multiplier

Obs. Value

1-1 Total mid-year population - Population totale au milieu de l'année

Thousands - milliers

Country - Pays	1980	1985	1990	1995	1999	2000	2001	2002	2003
Angola.....	6993	8754	9194	11072	12692	13134	13533	13942	14366
Botswana.....	906	1083	1276	1487	1529	1541	1549	1552	1565
Lesotho.....	1339	1538	1792	2050	2037	2035	2050	2065	2080
Malawi.....	6183	7340	9667	11129	11270	11308	11554	11806	12064
Mauritius - Maurice.....	966	1020	1057	1117	1151	1161	1169	1178	1187
Mozambique.....	12095	13711	14187	16004	17808	18292	18616	18946	19283
Namibia - Namibie.....	1030	1518	1349	1540	1711	1757	1787	1817	1848
South Africa	29170	33043	37066	41465	42902	43309	43634	43966	44306
Swaziland.....	560	664	744	855	910	925	933	942	950
Zambia - Zambie.....	5738	7006	8152	9456	10218	10421	10639	10683	11092
Zimbabwe.....	7126	8292	9903	11261	12333	12627	12843	13065	13292
Southern Africa, Total -									
Afrique de australe, totale.....	72106	83969	94387	107436	114561	116510	118305	119962	122033

SDMX Concept Scheme

- “Set of Concepts that are used in a Data Structure Definition or Metadata Structure Definition.”*
- Concept scheme places concepts into a maintainable unit.

Concept name	Concept ID
Indicator	INDICATOR
Reference area	REF_AREA
Time period	TIME_PERIOD
Unit multiplier	UNIT_MULT
Observation value	OBS_VALUE

Dimension

- Which of the concepts are used to identify an observation?
 - Indicator
 - Reference area
 - Time Period
- When all 3 are known, we can unambiguously locate an observation in the table.
- These are called **dimensions**.
 - A dimension is similar in meaning to a database table's primary key field.

Special Dimensions

- **TIME** dimension provides the period of time to which the observation relates. If a DSD describes time series data, it must have one TIME dimension.
 - **REFERENCE AREA** dimension describes the geographic location to which the observation refers (e.g., country, region, city, ...)
- **Everything happens at a specific moment (or period) in time**
- **Everything happens somewhere**

Attribute

- In our example, **Unit Multiplier** represents additional information about observations.
- This concept is not used to identify a series or observation.
- Such concepts in are called **attributes**.
 - Not to be confused with XML attributes!
 - Similar to a database table's non-primary key fields.

Primary Measure

- Observation Value represents a concept that describes the actual values being transmitted.
- In SDMX, such a concept is called **Primary Measure**.
- Primary Measure is usually represented by concept with ID **OBS_VALUE**.

Dimension or Attribute?

- Choosing the role of a concept has profound implications on the structure of data.
- Concepts that identify data, should be made dimensions. Concepts that provide additional information about data, should be made attributes.
- If a concept is a dimension, it is possible to have time series that are different only in the value of this concept.
 - E.g. if Unit of Measure is a dimension, it is possible to have separate series for "T" and "T/HA" or, more controversially, "KG" and "T"

Dimension or Attribute? (2)

Cambodia		
Fixed and Mobile telephone subscriptions	2013	20.6 million
Fixed and Mobile telephone subscriptions	2012	19.7 million
Fixed and Mobile telephone subscriptions	2013	140.9 per 100 pop.

Unit of measure as a dimension...

<u>Ref.Area</u>	<u>Indicator</u>	<u>Time Period</u>	<u>Unit</u>	<u>Unit Mult.</u>	<u>Obs. Value</u>
Cambodia	Fixed and Mobile telephone subscriptions	2013	Number	Millions	20.6
Cambodia	Fixed and Mobile telephone subscriptions	2012	Number	Millions	19.7
Cambodia	Fixed and Mobile telephone subscriptions	2013	Per 100 pop.	Units	140.9

22

Dimension or Attribute? (3)

Unit of measure as an attribute...

Violation!

<u>Ref.Area</u>	<u>Indicator</u>	<u>Time Period</u>	Unit	Unit Mult.	Obs. Value
Cambodia	Fixed and Mobile telephone subscriptions	2013	Number	Millions	20.6
Cambodia	Fixed and Mobile telephone subscriptions	2012	Number	Millions	19.7
Cambodia	Fixed and Mobile telephone subscriptions	2013	Per 100 pop.	Units	140.9

- The dataset above is invalid: duplicate observation
- The two values above are only different in their attributes

Dimension or Attribute? (4)

Unit of measure as an attribute...



<u>Ref.Area</u>	<u>Indicator</u>	<u>Time Period</u>	Unit	Unit Mult.	Obs. Value
Cambodia	Fixed and Mobile telephone subscriptions	2013	Number	Millions	20.6
Cambodia	Fixed and Mobile telephone subscriptions	2012	Number	Millions	19.7
Cambodia	Fixed and Mobile telephone subscriptions per 100 population	2013	Per 100 pop.	Units	140.9

- Now there is no violation because every row has a unique key
- The Unit concept is still useful

Attribute attachment

- In SDMX 2.0, attributes can be attached at observation, time series, group, or dataset level.
- In SDMX 2.1, attributes can be attached at observation, dimension(s), group, or dataset.
 - When attribute is attached to all dimensions except time, it is effectively attached to time series
- For practical purposes attributes are often attached at observation or time series.

Data model so far...

Concept	ID	Role	Attachment
Indicator	INDICATOR	Dimension	
Reference area	REF_AREA	Dimension	
Time period	TIME_PERIOD	Dimension	
Unit multiplier	UNIT_MULT	Attribute	Time series
Observation value	OBS_VALUE	P.Measure	

Exercise 1: Identifying concepts

- Identify concepts in the table
- Mark each concept as:
 - Dimension
 - Primary Measure
 - Attribute
- Identify the Time Dimension (Reference Period)
- Identify the Reference Area Dimension

Representation

- DSD defines a range of valid values for each concept.
- When data are transferred, each of its descriptor concepts must have valid values.
- A concept can be
 - Coded
 - Un-coded with format
 - Un-coded free text

Code

- "A language-independent set of letters, numbers or symbols that represent a concept whose meaning is described in a natural language."
- A sequence of characters that can be associated with a descriptions in any number of languages.
 - Descriptions can be updated without disrupting mappings or other components of data exchange.

Code List

- “A predefined list from which some statistical coded concepts take their values.”
- A code list is a collection of codes maintained as a unit.
- A code list enumerates all possible values for a concept or set of concepts
 - Sex code list
 - Country code list
 - Indicator code list, etc

Code List: Some Examples

CL_SERIES	
Code	Description
SI_POV_DAY1	Population below international poverty line (1.1.1)
SI_POV_EMP1	Employed population below international poverty line (1.1.1)
SI_POV_NAHC	Population below national poverty line (1.2.1)
SI_COV_BENFTS	Population covered by at least one social protection floor/system (1.3.1)
SI_COV_CHLD	Children covered by social protection (1.3.1)
SI_COV_DISAB	Population with severe disabilities collecting disability social protection benefits (1.3.1)
SI_COV_LMKT	Population covered by labour market programs (1.3.1)
SI_COV_MATNL	Mothers receiving maternity benefits and benefits for newborns (1.3.1)
SI_COV_PENSN	Population above retirement age receiving a pension (1.3.1)

CL_EDUCATION_LEV		
Code	Description (EN)	Description (FR)
_T	Total or no breakdown by education level	Total ou aucune ventilation par niveau de s
ISCED11_0	Early childhood education	Education de la petite enfance
ISCED11_01	Early childhood educational development	Développement éducatif de la petite enfance
ISCED11_02	Pre-primary education	Enseignement préprimaire
ISCED11_1	Primary education	Enseignement primaire
ISCED11_10	Primary education	Enseignement primaire

CL_REF_AREA	
Code	Description
1	World
2	Africa (M49)
4	Afghanistan
5	South America (M49)
8	Albania
9	Oceania (M49)
10	Antarctica
11	Western Africa (M49)
12	Algeria

SDMX Concepts and Code lists

- Code lists provide a representation for concepts, in terms of Codes.
- Codes are language-independent and may include descriptions in multiple languages.
- Code lists must be harmonized among all data providers that will be involved in exchange.

Un-coded Concepts

- Can be free-text: Any valid text can be used as a value for the concept.
 - Footnote
- Can have their format specified
 - Postal code: 5 digits
 - Last update: date/time

Representation of concepts in SDMX

- **Dimensions** must be either coded or have their format specified.
 - Free text is not allowed.
- **Attributes** can be coded or un-coded; format may optionally be specified.

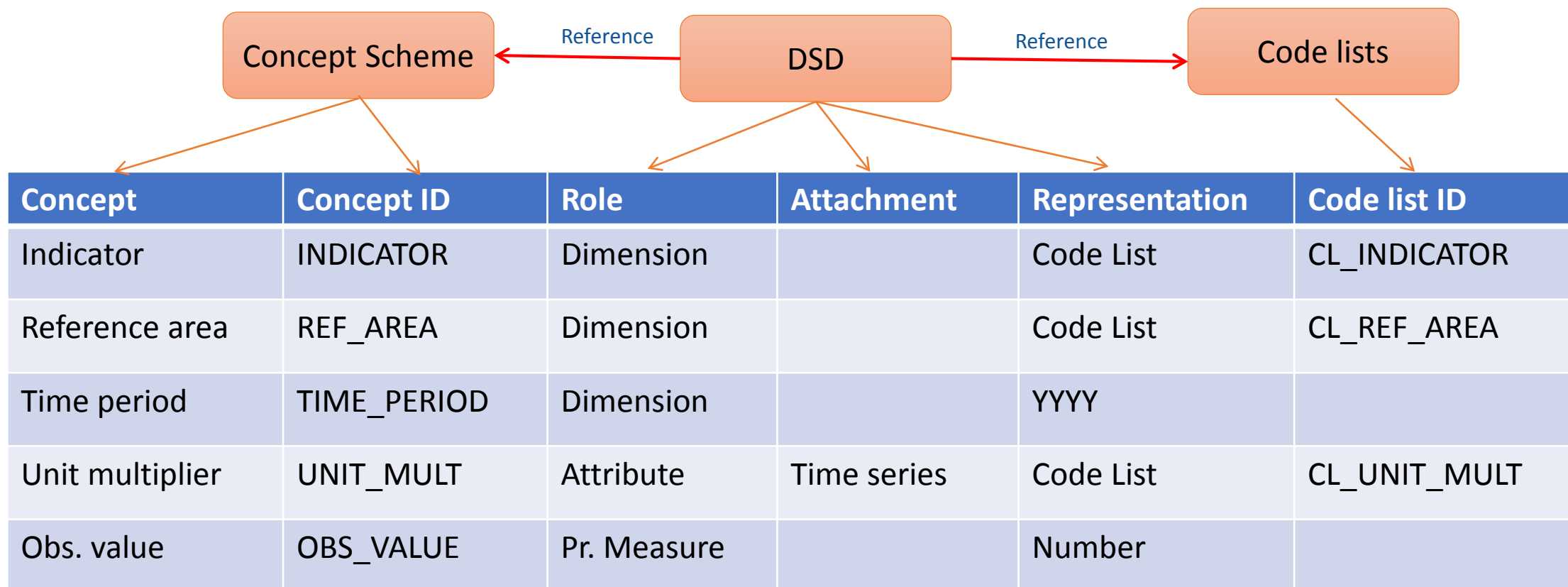
Data model so far...

Concept	ID	Role	Attachment	Representation
Indicator	INDICATOR	Dimension		CL_INDICATOR
Reference area	REF_AREA	Dimension		CL_REF_AREA
Time period	TIME_PERIOD	Dimension		Date/time (YYYY)
Unit multiplier	UNIT_MULT	Attribute	Time series	CL_UNIT_MULT
Observation value	OBS_VALUE	Pr. Measure		Floating point number

Exercise 2: Representation

- Working with your model, determine representation for each concept
 - Coded, formatted, free-text
- Develop code lists and formats for your concepts
 - Choose any approach for your codes and use it consistently

Data Structure Definition: summary



Importance of Data Model

- Data model, represented by DSD, defines what data can be encoded and transmitted.
- Flaws in a DSD may have significant adverse impact on data exchange
 - Missing concepts
 - Incorrect role of concepts
 - Un-optimized model

Data Structure Definition: Design Considerations

- Parsimony
 - No redundant dimensions
 - Attributes attached at the highest possible level
- Simplicity
 - “Mixed dimensions” are used to minimize the number of dimensions
 - Can help avoid invalid combinations of key values
 - Should be used with caution
 - Opposite of “purity”

Data Structure Definition: Design Considerations (2)

- Unambiguousness
 - Data must retain meaning outside usual context
 - Do you supply country code with your data?
- Density
 - Model should be such that data could be supplied for most or all of possible combinations of key values
 - Related to simplicity
- Orthogonality
 - Meaning of the value of concepts should be independent of each other
 - Helps avoid ambiguity

DSD Design Tradeoffs: Simplicity vs Purity

- A *simple* model may increase maintenance costs
 - Codes frequently need to be added
 - Difficult to map and consume
- A *pure* model may increase the number of errors due its lower *density*
 - Some combinations of key values are impossible in reality but valid from the DSD point of view
- Splitting the *pure* model into multiple DSDs to improve *density* may increase maintenance costs
 - Multiple DSDs and other artefacts need to be maintained

Dataset

- Organised collection of data defined by a Data Structure Definition (DSD)*
- A dataset is structured in accordance with *one* DSD
- Serves as a container for time-series or cross-sectional series in SDMX data messages.

Time Series

- A set of observations of a particular variable, taken at different points in time.
- Observations that belong to the same time series, differ in their time dimension.
 - All other dimension values are identical.
 - Observation-level attributes may differ across observations of the same time series.

Time Series: Demonstration

1.1 Proportion of population below \$1 (PPP) per day													
Series	1990	1992	1994	1996	1998	1999	2000	2002	2006	2007	2008	2009	2011
Rwanda													
MDG Population below \$1 (PPP) per day, percentage Last updated: 02 Jul 2012							74.6 ^{1,3}		72.1 ^{1,3}				63.2 ^{1,3}
State of Palestine													
MDG Population below \$1 (PPP) per day, percentage Last updated: 02 Jul 2012										0.4 ^{1,2,3}		0.0 ^{1,2,3}	
Thailand													
MDG Population below \$1 (PPP) per day, percentage Last updated: 02 Jul 2012	11.6 ^{1,3}	8.6 ^{1,3}	4.1 ^{1,3}	2.5 ^{1,3}	2.1 ^{1,3}	3.2 ^{1,3}	3.0 ^{1,3}	1.6 ^{1,3}	1.0 ^{1,3}		0.4 ^{1,3}	0.4 ^{1,3}	
1.2 Poverty gap ratio													
Series	1990	1992	1994	1996	1998	1999	2000	2002	2006	2007	2008	2009	2011
Rwanda													
MDG Poverty gap ratio at \$1 a day (PPP), percentage Last updated: 02 Jul 2012							36.9 ^{1,3}		34.8 ^{1,3}				26.6 ^{1,3}
State of Palestine													
MDG Poverty gap ratio at \$1 a day (PPP), percentage Last updated: 02 Jul 2012										0.1 ^{1,2,3}		0.0 ^{1,2,3}	
Thailand													
MDG Poverty gap ratio at \$1 a day (PPP), percentage Last updated: 02 Jul 2012	2.4 ^{1,3}	1.6 ^{1,3}	0.7 ^{1,3}	0.4 ^{1,3}	0.3 ^{1,3}	0.5 ^{1,3}	0.5 ^{1,3}	0.3 ^{1,3}	0.2 ^{1,3}		0.0 ^{1,3}	0.1 ^{1,3}	
Footnotes													
1 Based on nominal per capita consumption averages and distributions estimated from household survey data.													
2 Based on Purchasing Power Parity (PPP) dollars imputed using regression.													
3 Source: http://research.worldbank.org/PovcalNet/index.htm													

Non-Time Series Data (a.k.a. Cross-Sectional Data)

- A non-time dimension is chosen along which a set of observations is constructed.
 - E.g. for a survey or census the time is usually fixed and another dimension may be chosen to be reported at the observation level
- Used less frequently than time series representation

Time Series View vs Cross-Sectional View

2.1 Net enrolment ratio in primary education

	2009	2010	2011
Morocco			
Total net enrolment ratio in primary education, both sexes		94.1	96.2
Total net enrolment ratio in primary education, boys		95	96.8
Total net enrolment ratio in primary education, girls		93.3	95.6
State of Palestine			
Total net enrolment ratio in primary education, both sexes	88.2	89.2	
Total net enrolment ratio in primary education, boys	88.2	89.8	
Total net enrolment ratio in primary education, girls	88.2	88.5	
Uganda			
Total net enrolment ratio in primary education, both sexes	94.2	91	
Total net enrolment ratio in primary education, boys	93.1	89.7	
Total net enrolment ratio in primary education, girls	95.3	92.3	

- The Sex dimension was chosen as the cross-sectional measure.

- Note that Time is still applicable.

2.1 Net enrolment ratio in primary education 2010

	Total	Boys	Girls
Morocco	94.1	95	93.3
State of Palestine	89.2	89.8	88.5
Uganda	91	89.7	92.3

Keys in SDMX

- **Series key** uniquely identifies a series
 - In the case of time series, consists of all dimensions except **time**
- **Group key** uniquely identifies a group of time series
 - Consists of a subset of the series key

Exercise 3: Encoding a time series

- Working with your table, determine the total number of time series.
- For the first 5 time series, provide a valid value for each concept in its series key.

Structural and Reference Metadata

- Structural Metadata: Identifiers and Descriptors, e.g.
 - Data Structure Definition
 - Concept Scheme
 - Code
- Reference Metadata: Describes contents and quality of data, e.g.
 - Indicator definition
 - Comments and limitations

} What we
have
covered
so far

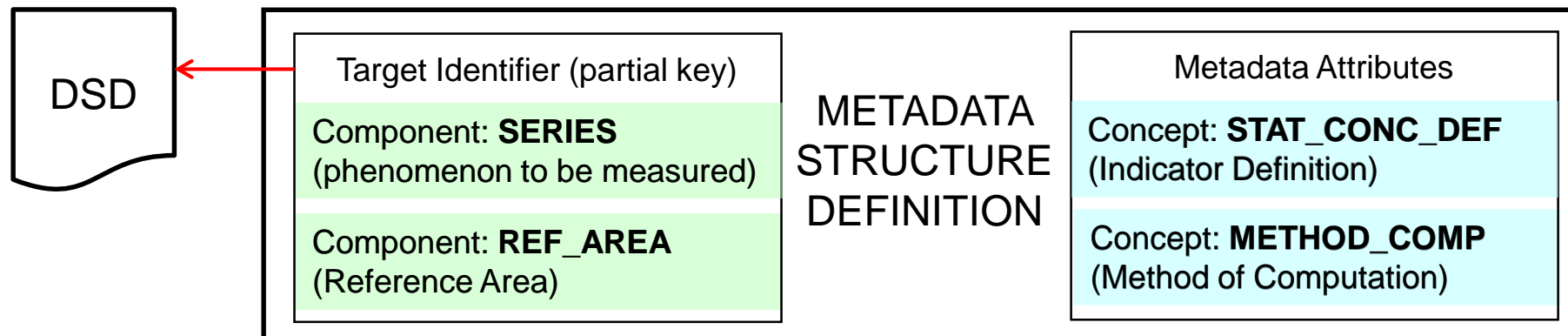
Reference Metadata in SDMX

- Can be stored or exchanged separately from the object it describes, but be linked to it
- Can be indexed and searched
- Reported according to a defined structure

Metadata Structure Definition (MSD)

- MSD Defines:
 - The object type to which reference metadata can be associated
 - E.g. DSD, Dimension, Partial Key.
 - The components comprising the object identifier of the target object
 - E.g. the draft SDG MSD allows metadata to be attached to each indicator for each country
- Concepts used to express metadata ("metadata attributes").
 - E.g. Indicator Definition, Quality Management

Metadata Structure Definition and Metadata Set: an example



METADATA SET
SERIES= SH_STA_BRTC (Births attended by skilled health personnel) REF_AREA= KH (Cambodia)
STAT_CONC_DEF ="It refers to the proportion of deliveries that were attended by skilled health personnel including physicians, medical assistants, midwives and nurses but excluding traditional birth attendants."
METHOD_COMP ="The number of women aged 15-49 with a live birth attended by skilled health personnel (doctors, nurses or midwives) during delivery is expressed as a percentage of women aged 15-49 with a live birth in the same period. "

Dataflow and Metadataflow

- Dataflow defines a “view” on a Data Structure Definition
 - Can be constrained to a subset of codes in any dimension
 - Can be categorized, i.e. can have *categories* attached
 - In its simplest form defines any data valid according to a DSD
- Similarly, Metadataflow defines a view on a Metadata Structure Definition.

Content Constraints

- Constraints can be used to define which combinations of codes are allowed
 - E.g. "When ***SERIES***= Proportion of Women in Commune Councils', ***SEX*** must be Female"
- Constraints can define more granular validation rules than a simple validation of codes
- Are often attached to the Dataflow but can also be attached to DSD, Provision Agreement, etc

Category and Category Scheme

- Category is a way of classifying data for reporting or dissemination
 - Subject matter-domains are commonly implemented as Categories, such as "Demographic Statistics", "Economic Statistics"
- Category Scheme groups Categories into a maintainable unit.

Dataflows - classification

Categories

Tourism

Capacity







Occupancy

Night_Spent

Arrival_of_residents

Occupancy_rate

Dataflows

 Number of establishments, bedrooms and bed-places - national - annual data (tour_cap_nat) 
 Number of establishments, bedrooms and bed-places by NUTS 3 regions - annual data (tour_cap_nuts3) 
 Bed-places (x1 000) (tour_cap_bed) 

 Nights spent in tourist accommodation establishments - national - monthly data (tour_occ_nim)  **Updated**
 Nights spent by non-residents in tourist accommodation establishments - 1990-2011 - world geographical breakdown - monthly data (tour_occ_ninmw) 
 Nights spent in tourist accommodation establishments - national - annual data (tour_occ_ninat) 
 Nights spent in tourist accommodation establishments by NUTS 2 regions - annual data (tour_occ_nin2) 
 Nights spent by non-residents in tourist accommodation establishments - world geographical breakdown - annual data (tour_occ_ninraw) 
 Nights spent (x1 000) (tour_occ_ni) 

SDMX Messages

- Any SDMX-related information is exchanged in the form of documents called *messages*. An SDMX message can be sent in a number of standard formats including XML, JSON, CSV
- There are several types of SDMX messages, each serving a particular purpose, e.g.
 - **Structure** message is used to transmit structural information such as DSD, MSD, Concept Scheme, etc.
 - **GenericData**, **StructureSpecificData**, and other messages are used to send data.
- SDMX messages in the XML format are referred to as SDMX-ML messages.



United Nations


Department of Economic and Social Affairs
Statistics

The SDG Data Structure Definition



SDG Data Structure Definition

- Developed by the Working Group on SDMX for SDG Indicators, established by IAEG-SDGs in April 2016
- First version officially released on 14 June 2019

 [SDG DSD Matrix Version 1.0](#)

 [Global SDG DSD v1.0](#)

 [Guidelines for the Global DSD for SDGs](#)

<https://unstats.un.org/sdgs/iaeg-sdgs/sdmx-working-group/>

SDG Data Structure Definition

- One single DSD is used for all SDG indicators
- Support for diverse indicators means not all dimensions are applicable in all cases
 - E.g. AGE is not applicable to indicator "Land area covered by forest"
 - Value **_T** (no breakdown) is used when an dimension is not applicable.

Dimension: Frequency (FREQ)

- "Indicates rate of recurrence at which observations occur (e.g. monthly, yearly, biannually, etc.)."
- By convention, SDGs DSD currently only supports annual frequency.
- Where the frequency is not annual (e.g. two-year average), detail should be provided in the TIME_DETAIL attribute.

Dimension: REPORTING_TYPE

- Used to distinguish between National, Regional, Global Reporting
- Countries to use value **N** (national reporting)
- Regional organizations to use value **R** (regional reporting)
- Custodian agencies to use value **G** (Global reporting)

Dimension: Series (SERIES)

- Used to represent “sub-indicators”
 - A single indicator can have multiple series
 - Not to be confused with SDMX time series (each series can have multiple time series, i.e., multiple disaggregation with observations organized over time)
- Example: Indicator 5.5.1, “Proportion of seats held by women in (a) national parliaments and (b) local governments” has 4 series:
 - SG_GEN_PARL Proportion of seats held by women in national parliaments
 - SG_GEN_PARLN Number of seats held by women in national parliaments
 - SG_GEN_PARLNT Number of seats in national parliaments
 - SG_GEN_LOCG Proportion of seats held by women in local governments

Dimension: Reference Area (REF_AREA)

- Country or geographic area to which the measured statistical phenomenon relates
- It is envisaged that countries will report national-level values but may wish to extend the code list with its sub-national areas for dissemination

Dimension: Sex (SEX)

- Gender condition: male or female. This dimension applies only if data can be disaggregated by sex.
- Use **_T** where not applicable
- For gender indicators must be set to **F** as applicable
 - E.g. for series *Proportion of seats held by women in national parliaments*

Dimension: Age (AGE)

- "Age - or age range - of the individuals the observation refers to."
- Use **_T** where not applicable

Dimension: Urban/Rural location (URBANISATION)

- Has 3 codes
 - _T (Total)
 - _U (Urban)
 - _R (Rural)
- Use _T where not applicable

Dimension: INCOME_WEALTH_QUANTILE

- Used for disaggregating the data by income or wealth quintile of the population
- In the future can be extended to cover decile, percentile, etc
- Use _T where not applicable

Dimension: Education Level (EDUCATION_LEV)

- “Highest level of an educational programme the person has successfully completed.”
- Supports top categories of ISCED11 and ISCED97, as well as custom SDG codes
- Use _T where not applicable

Dimension: OCCUPATION

- “Job or position held by an individual who performs a set of tasks and duties.”
- Supports top categories of ISCO-08, ISCO-98, ISCO-68
- Use _T where not applicable

Dimension: Disability Status (DISABILITY STATUS)

- Used to break down SDG indicators by disability
- At the moment, only used to distinguish between persons with a disability, and persons without a disability
- Use _T where not applicable

Dimension: Economic Activity (ACTIVITY)

- “High-level grouping of economic activities based on the types of goods and services produced.”
- Consists of top-level ISIC categories.
- Use **_T** where not applicable.

Dimension: Product Type (PRODUCT)

- Product or commodity code
- Combines SDG-specific entries from several classifications including CPC, Material Flows, and non-standard
- Use **_T** where not applicable

Dimension: Custom Breakdown (CUST_BREAKDOWN)

- Special dimension introduced to facilitate non-standard breakdowns, primarily in national context
- At the moment empty but in the future will be populated with generic codes (e.g. CODE1, CODE2, etc), to which data providers will assign meaning in their own context
- Used in conjunction with attribute CUST_BREAKDOWN_LB, which transmits description of the custom code.
- Use **_T** where not applicable

Dimension: COMPOSITE_BREAKDOWN

- Mixed dimension: represents several merged code lists
 - E.g. by International Organizations, Hazard Type etc
- Used for breakdowns that are only used in 1 or 2 indicators, in order to avoid creating too many dimensions
- Use _T where not applicable

Time Dimension: TIME_PERIOD

- The observation corresponds to a specific point in time ... or a period...”
- The convention for SDGs is to always provide a four-digit year in the TIME_PERIOD concept. Further info must be placed in TIME_DETAIL, and structured period information in TIME_COVERAGE.

Primary Measure: Observation value (OBS_VALUE)

- Used to convey the value of a variable at a period of time
- Should be a floating-point number

Attribute: Observation Status (OBS_STATUS)

- "Information on the quality of a value or an unusual or missing value"
 - E.g. can be used to indicate a break in series
- Mandatory observation-level attribute

Attribute: Unit Multiplier (UNIT_MULT)

- Exponent in base 10 specified so that multiplying the observation numeric values by **$10^{\text{UNIT_MULT}}$** gives a value expressed in the unit of measure
- If the observation value is in millions, unit multiplier is 6; if in billions, 9, and so on. Where the number is simple units, use 0.
- Mandatory observation-level attribute

Attribute: Unit of Measure (UNIT_MEASURE)

- Unit in which the data values are expressed
- It may not be obvious which is the correct unit in some cases. Coding guidelines are available and will be further developed.
- Mandatory time series-level attribute

Attribute: Time Period Details (TIME_DETAIL)

- “When TIME_PERIOD refers to a date range, this attribute is used to provide metadata on the actual range the observation refers to (e.g. for period ‘2001-2003’ TIME_PERIOD would be 2002 but the actual dates --2001-2003-- would be expressed here).”
- Optional observation-level free-text attribute

Attribute: TIME_COVERAGE

- ISO8601 representation of the actual time interval to which the observation refers
- While TIME_PERIOD should always be expressed as a year, and TIME_DETAIL is free-text with additional information, TIME_COVERAGE can optionally be used to provide the exact interval in a structured format
- Optional observation-level attribute.

Attribute: Base Period (BASE_PER)

- Period of time used as the base of an index number, or to which a constant series refers
- Where a base period applies, it is expected to always be set to a year
- Typically, used for constant prices, as in "2005 USD dollar"
- Optional observation-level attribute.

Attribute: Nature of data points (NATURE)

- Information on the production and dissemination of the data
- Expresses whether a data point has been produced and disseminated by the country, estimated by international agencies, etc.
- Normally set to C (Country Data) in national reporting
- Optional observation-level attribute

Attribute: Source details (SOURCE_DETAIL)

- Provides additional textual information on the data source, e.g. a specific survey that was used to generate the indicator.
- Optional observation-level free-text attribute.

Attributes: UPPER_BOUND and LOWER_BOUND

- Where the observation value represents a point estimate, can be used to convey the Upper and Lower bounds
 - In SDG DSD, separate series codes had to be created for upper and lower bounds
- Optional observation-level attributes

Attributes: Footnotes (COMMENT_OBS and COMMENT_TS)

- “Additional information on specific aspects of each observation, such as how the observation was computed/estimated or details that could affect the comparability of this data point with others in a time series.”
- Attribute COMMENT_OBS is used for observation-level footnotes, and COMMENT_TS for time series-level footnotes. Both are optional.

Attribute: GEO_INFO_URL

- Provides web address of a geoinformation file. Used in conjunction with attribute GEO_INFO_TYPE.
- Optional time series-level attribute.

Attribute: GEO_INFO_TYPE

- Specifies type of geoinformation file provided in attribute GEO_INFO_URL.
- Optional time series-level attribute.

SDG DSD: Mappings

- Due to its support for heterogeneous indicators, it's not always obvious which values should be used in some dimensions
- What should be SEX in indicator "Births attended by skilled personnel":
 - Not Applicable? Total? Female?

SDG DSD: Mappings

- Inconsistent mappings lead to duplications and other anomalies
- Coding guidelines will be developed and enforced through content constraints
- The use of a single code for no breakdown (e.g. for Total and Not Applicable) simplifies the mappings.

Exercise 1: Identifying concepts

- Identify concepts in the table
- Mark each concept as:
 - Dimension
 - Primary Measure
 - Attribute
- Identify the Time Dimension (Reference Period)
- Identify the Reference Area Dimension

Exercise 2: Representation

- Working with your model, determine representation for each concept
 - Coded, formatted, free-text
- Develop code lists and formats for your concepts
 - Choose any approach for your codes and use it consistently

Exercise 3:

Encoding a time series

- Working with your table, determine the total number of time series.
- For the first 5 time series, provide a valid value for each concept in its series key.



**United
Nations**

DESA
Statistics Division

Thank you.