



# Data quality and verification

---

Amani Al-Saidi

Training Workshops on Energy Statistics: Energy Balances

December 13-14, 2018, Beirut





Balance checks



Time series analysis



Efficiencies



Data validation in practice

# What is data quality?

*Data quality* is most commonly assessed in terms of its “**fitness for use**”.

How well do the statistical outputs meet user needs?

Relevance	Statistics meet the <b>needs of the users</b>
Accuracy and reliability	Statistics accurately and reliably portray <b>reality</b>
Timeliness and punctuality	Statistics are released in a <b>timely</b> and <b>punctual</b> manner
Accessibility and clarity	Statistics are presented in a clear and <b>understandable</b> form, released in a suitable and convenient manner, available and accessible on an impartial basis with supporting metadata and <b>guidance</b> .
Coherence and comparability	Statistics are <b>consistent</b> internally, over time and <b>comparable between regions and countries</b> .

For more information please see:

[Fundamental Principles of Official Statistics](#) and [International Recommendation for Energy Statistics \(IRES\)](#)



# Balance checks

---

Internal consistency

Statistical differences

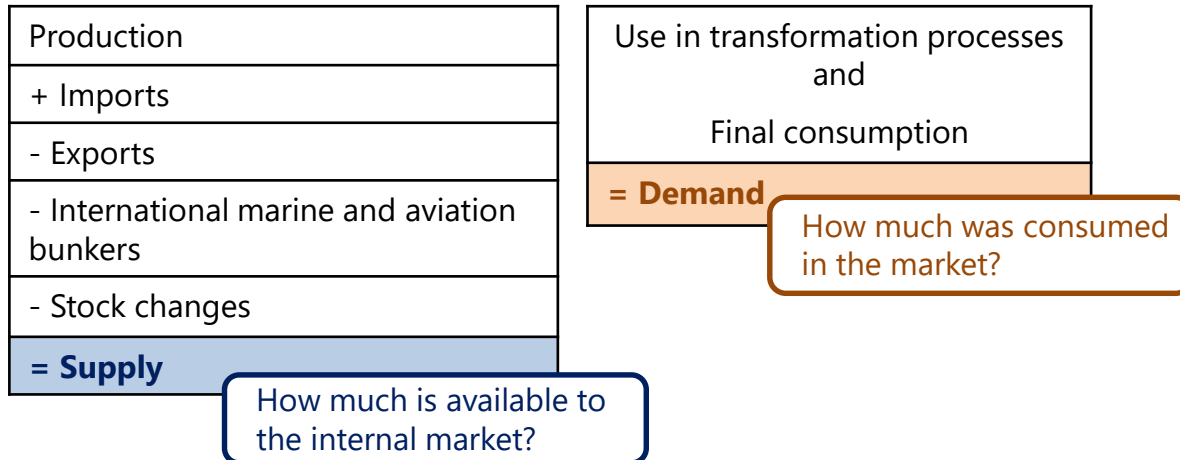
Own use and losses

Calorific values

- Some things to look out for...
  - If one flow or product is supposed to be the sum of two others , is it?
  - Does the detailed breakdown match the total (e.g. for trade or demand data)?
  - Are stock changes correctly calculated?
  - Is the data set complete?
  - Decrease in use of one fuel coupled with increase in use of another one?
  - Is the story logical?
- It is good practice to analyze the data processes to identify possible points of failure and introduce checks.

# Balance checks - What are statistical differences?

- When collecting energy data we try to get a picture of supply and demand in a market.
  - Economic theory tells us that when a market is in equilibrium supply = demand.
- In practice, measuring supply and demand comes with many challenges.
- Statistical differences measure the level of mismatch between the collected supply and demand data.



- High statistical differences can be a red flag (often when more than 5% of deliveries to the market).
- But... a statistical difference of 0 is also suspicious!
- Calculating the **ratio of statistical difference to total primary energy supply (TPES)** can help us evaluate the size of the problem.

### *What can we do?*

Check the completeness and coverage of the data

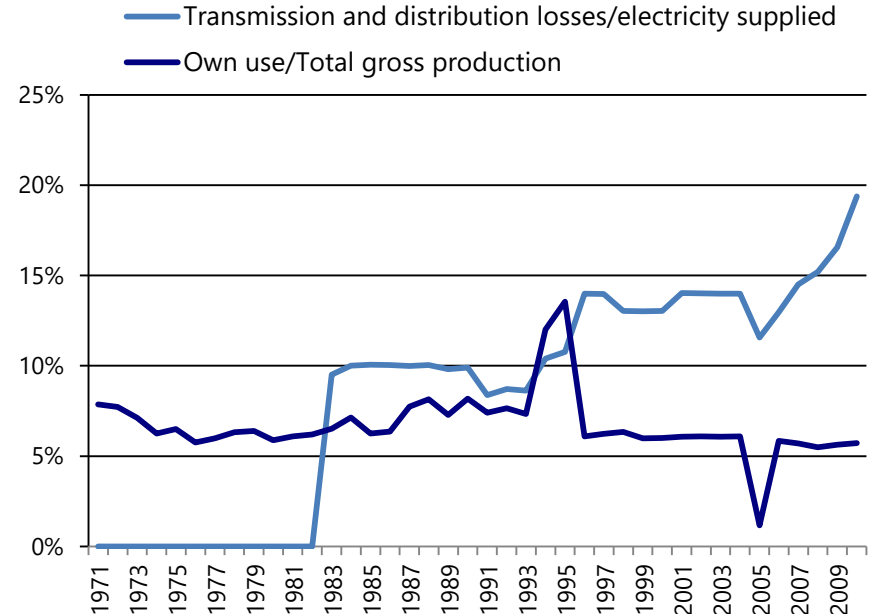
Review definitions and methods of data sources

Check calorific values for supply and demand

Are there any new trends missing? (renewables, electric vehicles...)

## Case study: Electricity and heat

- **Own use** of electricity by power plants (*main electricity plants only*) as a proportion is not expected to vary significantly year-on-year ( $\leq 3\%$ )
- **Transmission and distribution losses** should technically stay **below 10%**.
  - Bad maintenance and theft can lead to higher transmission and distribution losses (~ 30%)





- Calorific values are essential to building an energy balance and for converting to/from energy units.
- For every energy product there is an expected range for the calorific values.
  - The collected Net Calorific Values (NCVs) should be compared against default NCV
  - And also with historical NCVs for the product as calorific values can be region specific (except for natural gas as it is mostly methane)
  - Check if NCVs used for consumption is in line with the one used for supply.

Fuel	Expected calorific value (kJ/kg, MJ/ton)			GCV estimation
Coking coal	25000	-	33000	≈ NCV + 5%
Anthracite	22000	-	29000	≈ NCV + 5%
Other bituminous coal	22000	-	29000	≈ NCV + 5%
Sub-bituminous coal	16000	-	24000	≈ NCV + 5%
Lignite	5000	-	18000	≈ NCV + 5%
Peat	7000	-	13000	≈ NCV + 5%
Oil Shale	2500	-	12000	≈ NCV + 5%

*Expected GCV (kJ/m<sup>3</sup>)  
range for natural gas*

**30 000 – 45 000**



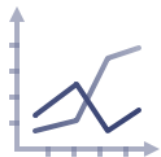
*Different qualities of bituminous coal around the world*



28 201 kJ/kg



18 581 kJ/kg



# Time series analysis

---

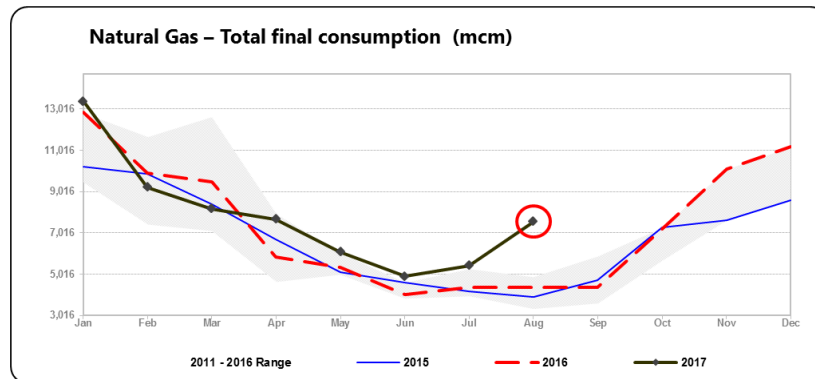
Growth rates

Time series breaks

Outliers

$$\text{Growth Rate (\%)} = \frac{\text{Current period} - \text{Previous period}}{\text{Previous period}} \times 100$$

- When seasonality is strong or we are dealing with monthly data
  - It is often more useful to compare the current month with the same month of the previous year.



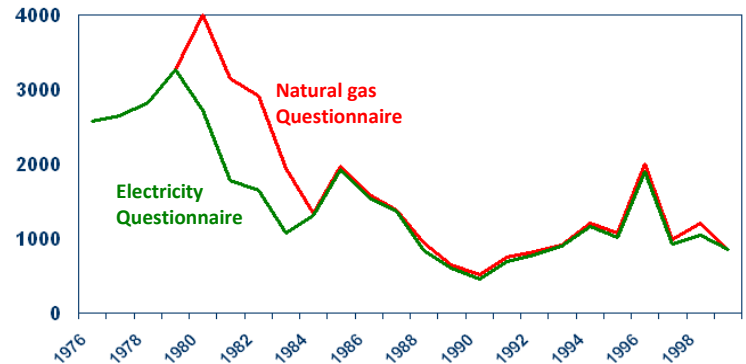
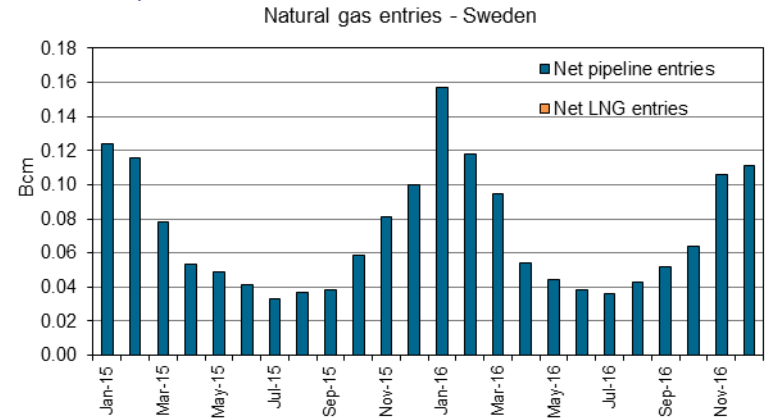
**Production  
Consumption**

*Less useful for*

**Stocks  
Trade**

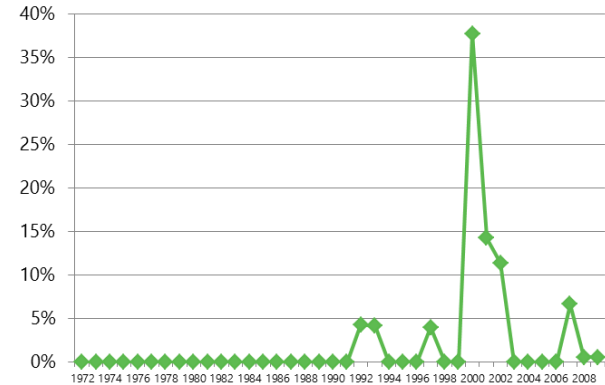
- Does the data set match expectations for trends? (e.g. known maintenance, weather disruptions, seasonality patterns, changes in policy...).
- Certain patterns may be fuel specific:
  - Natural gas demand and storage is linked to weather patterns.
  - Oil and gas fields follow cyclical maintenance patterns.
- Does the data match similar data reported in other questionnaires?
  - E.g. does inputs for electricity and heat match in the fuel and electricity data sets?

**Example.** Sweden imports more natural gas in the winter months to cope with increasing demand that comes from colder temperatures.



- Although breaks can happen when coverage improves or methodology changes in general one should aim for a **consistent time series**.
- **Unexplained** breaks or big changes in trends should be avoided.
- **Breaks in stocks** are rare and should be investigated.
- Certain flows are generally consistent over time, thus sudden increases or decreases should raise a red flag.

Biofuels and waste – TPES growth rate



Barrels per day Refinery Intake



A graphic consisting of two overlapping arrows: a blue arrow pointing upwards and a red arrow pointing downwards, positioned to the left of the title.

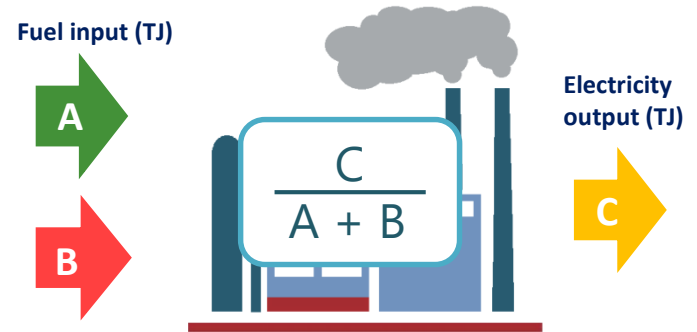
# Efficiencies

---

- For every energy transformation process there is an **input** and an **output**.
  - There is a relationship between the inputs and outputs which will depend on the technology used.
- Both inputs and outputs need to be in:
  - The same energy unit.
  - The same energy form (net or gross basis)
- This efficiency then needs to be checked against the expected ranges for that process and technology.
  - An efficiency larger than 100% is implausible.

$$\text{Efficiency} = \frac{\text{Total "useful" output (net energy units)}}{\text{Total input (net energy units)}}$$

**Example.** Fuel-electricity transformation



- Calculating and checking efficiencies allows us to check the **consistency** of our data set.
- Efficiencies are expected to be within a certain range depending on the **type of plant** and the **fuel used**.
- The ranges are wide because efficiency can vary widely depending on the **state of technology** (old v. new).
- Efficiencies **do not vary** significantly from year to year as technology upgrades are time and capital intensive.

Type of plant	Expected Efficiency range
Electricity only plants	10 - 40%
CHP plants	30-80%
Heat only plants	40-100%
Refineries	95-100%
Blast furnaces	35-45%
Coke ovens	67-100%* *(coke oven coke + coke oven gas)
Patent fuel plants	90-100%
BKB	85-100%
Gas works	67-100%* *(gas works gas + gas coke)
Charcoal	25-55%

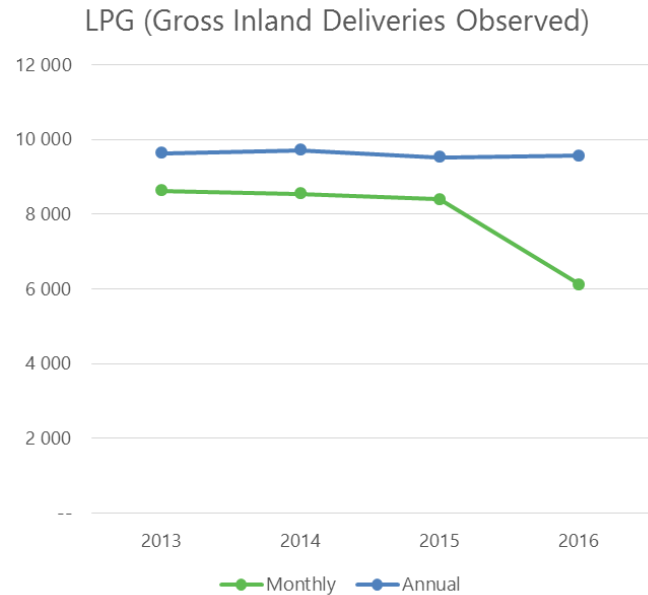


Potential issue	What can be done...
Incomplete <b>input</b> or <b>output</b> data (measurement errors, different source institution...)	Check the correct <b>definitions</b> are being used. Missing cells in energy balance, check <b>logical flows</b> for the fuel.
<b>Heat sold output</b> for CHP and heat only plants is not being correctly accounted for.	Review <b>methodology for heat output</b> estimation.
<b>Conversion</b> to energy units of either the input or output.	Are all the flows being converted on a <b>net basis</b> (NCV)? Is the <b>calorific value</b> data within the expected ranges?
Incorrect accounting of industry <b>own use</b> and <b>output sold</b> .	Check <b>boundaries</b> of statistical reporting and that the correct definitions are being used.

- Good quality monthly data can be a tool for data validation as well as key for energy security purposes.
  - Does the aggregate of the monthly/quarterly data make sense with the annual data?

Some things to look out for...

- Revisions
- Different definitions or coverage due to timing constraints.



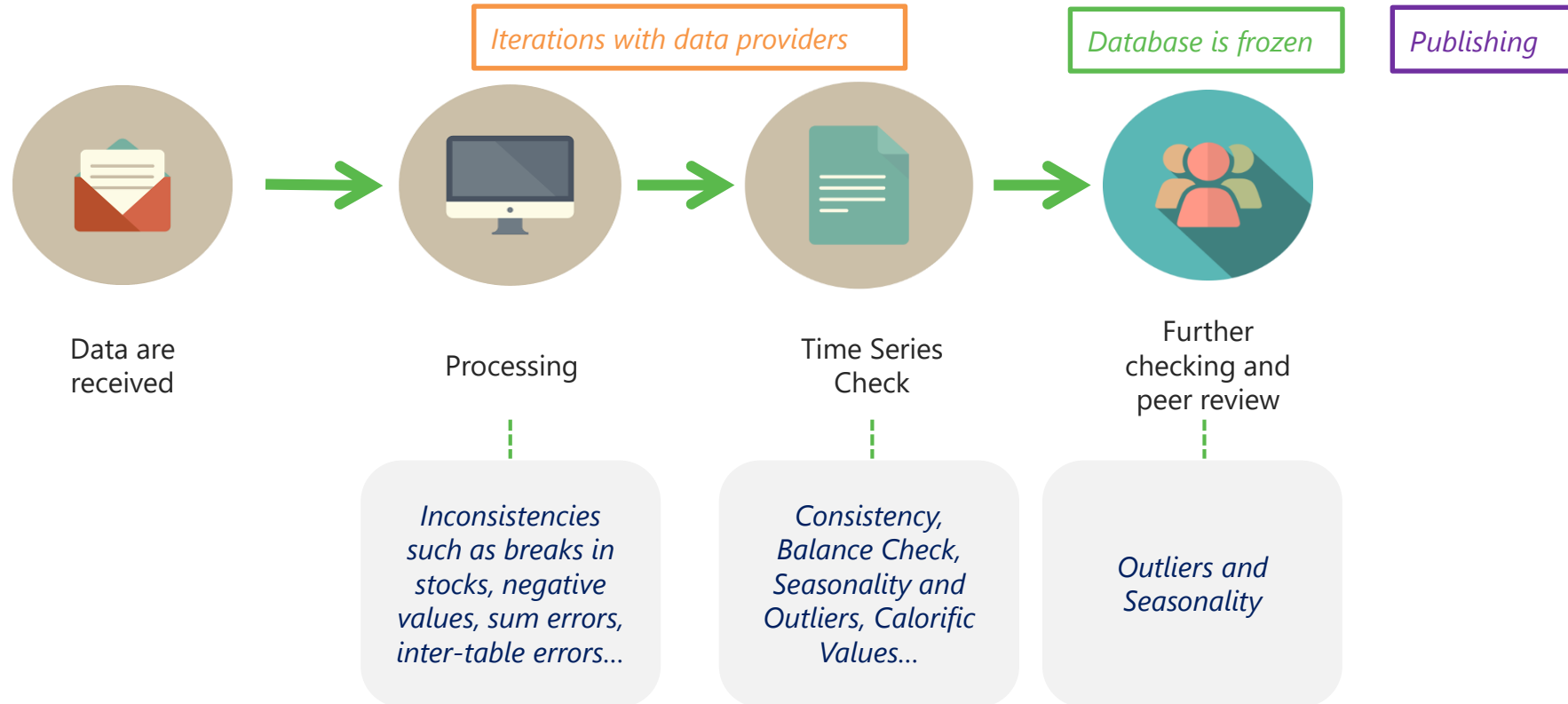


# Data validation in practice

---

# Data validation processes - Approaches

*An overview of how data is processed and validated at the IEA*



# Practical exercise – Checking the energy balance



- In groups discuss the different scenarios for checking the energy balance sheet.
- As a group, go through the questions raised in the files for each scenario, thinking about:
  - Checks that could be applied for checking the data when received?
  - What are some of the systems that are in place in your countries?

HOW TO CHECK AN ENERGY BALANCE?										Select language / Выбор языка: English	
Scenario 1: High statistical difference											
Statisland											
2015											
Thousand tonnes of oil equivalent (ktoe)											
SUPPLY AND CONSUMPTION	Coal	Crude oil	Oil products	Natural gas	Nuclear	Hydro	Geotherm./Solar/etc.	Biofuels/Waste	Electricity	Heat	Total
Production	15000	1000	-	-	-	-	-	-	-	-	16000
Imports	-	-	10000	5000	-	-	-	-	2000	-	17000
Exports	-	-	-	-	-	-	-	-	-	-	-
Int. marine bunkers	-	-	-	-	-	-	-	-	-	-	-
Int. aviation bunkers	-	-	-	-	-	-	-	-	-	-	-
Stock changes	-	-	-	250	-	-	-	-	-	-	250
<b>TOTAL</b>	<b>16000</b>	<b>1000</b>	<b>10000</b>	<b>6250</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>2000</b>	<b>-</b>	<b>33250</b>
Transfers	-	-	-	-	-	-	-	-	-	-	-
Statistical differences	2500	5	5000	-7000	-	-	-	-	700	-	1205
Electricity plants	-7500	-	-	-	-	-	-	-	2500	-	-5000
CHP plants	-	-	-	-5000	-	-	-	-	1300	2700	-1000
Heat plants	-	-	-	-	-	-	-	-	-	-	-
Blast furnaces	-	-	-	-	-	-	-	-	-	-	-
Gas works	-	-	-	-	-	-	-	-	-	-	-
Coal gas. Act/BK/IB plants	-	-	-	-	-	-	-	-	-	-	-
Oil refineries	-	-	-955	950	-	-	-	-	-	-	-45
Petrochemical plants	-	-	-	-	-	-	-	-	-	-	-
Liquefaction plants	-	-	-	-	-	-	-	-	-	-	-
Other transformation	-	-	-	-	-	-	-	-	-	-	-
Energy industry own use	-	-	-	-	-	-	-	-	-	-	-
Losses	-	-	-	-	-	-	-	-	-	-	-
<b>TFC</b>	<b>6000</b>	<b>-</b>	<b>8850</b>	<b>7250</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>6100</b>	<b>2700</b>	<b>26000</b>
<b>INDUSTRY</b>	<b>6000</b>	<b>-</b>	<b>3850</b>	<b>3750</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>3616</b>	<b>1716</b>	<b>17926</b>
Iron and steel	2500	-	1700	-	-	-	-	-	1250	627	6077
Chemical and petrochemical	-	-	1000	2500	-	-	-	-	450	655	4705
Non-ferrous metals	1500	-	700	-	-	-	-	-	350	291	2841
Non-metallic minerals	-	-	250	-	-	-	-	-	63	100	413
Transport equipment	-	-	100	500	-	-	-	-	32	55	688
Machinery	-	-	-	-	-	-	-	-	250	18	268
Mining and quarrying	1000	-	-	500	-	-	-	-	644	32	2176
Food and tobacco	-	-	-	-	-	-	-	-	250	16	266
Paper, pulp and printing	-	-	-	250	-	-	-	-	76	8	334
Wood and wood products	-	-	-	-	-	-	-	-	62	4	66
Construction	-	-	-	-	-	-	-	-	64	3	67
Textile and leather	-	-	-	-	-	-	-	-	26	1	27
Non-specified	-	-	-	-	-	-	-	-	8	-	8
<b>TRANSPORT</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>408</b>	<b>-</b>	<b>408</b>
Domestic aviation	-	-	-	-	-	-	-	-	-	-	-
Road	-	-	-	-	-	-	-	-	50	-	50
Rail	-	-	-	-	-	-	-	-	356	-	356
Pipeline transport	-	-	-	-	-	-	-	-	-	-	-
Domestic navigation	-	-	-	-	-	-	-	-	-	-	-
Non-specified	-	-	-	-	-	-	-	-	-	-	-
<b>OTHER</b>	<b>-</b>	<b>-</b>	<b>2000</b>	<b>2500</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>1178</b>	<b>890</b>	<b>8988</b>
Residential	-	-	1000	2500	-	-	-	-	782	740	5022
Comm. and public services	-	-	-	-	-	-	-	-	397	250	647
Agriculture/forestry	-	-	1000	-	-	-	-	-	-	-	1000
Fishing	-	-	-	-	-	-	-	-	-	-	-
Non-specified	-	-	-	-	-	-	-	-	-	-	-
<b>NON-ENERGY USE</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>1000</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>1000</b>
In industry/transf.energy	-	-	-	1000	-	-	-	-	-	-	1000
of which: chem./petrochem.	-	-	-	-	-	-	-	-	-	-	-
In transport	-	-	-	-	-	-	-	-	-	-	-
In other	-	-	-	-	-	-	-	-	-	-	-
<b>Electricity and heat output</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>16116</b>	<b>-</b>	<b>16116</b>
Electr. Generated - GWh	29070	-	-	-	-	-	-	-	-	-	29070
Electricity plants	29070	-	-	-	-	-	-	-	-	-	29070
CHP plants	-	-	-	16116	-	-	-	-	-	-	16116
Heat generated - TJ	-	-	-	113066	-	-	-	-	-	-	113066
CHP plants	-	-	-	113066	-	-	-	-	-	-	113066
Heat plants	-	-	-	-	-	-	-	-	-	-	-



[www.iea.org/statistics](http://www.iea.org/statistics)

