# Business intelligence

## Tutorial 3 – Pentaho Data Integration
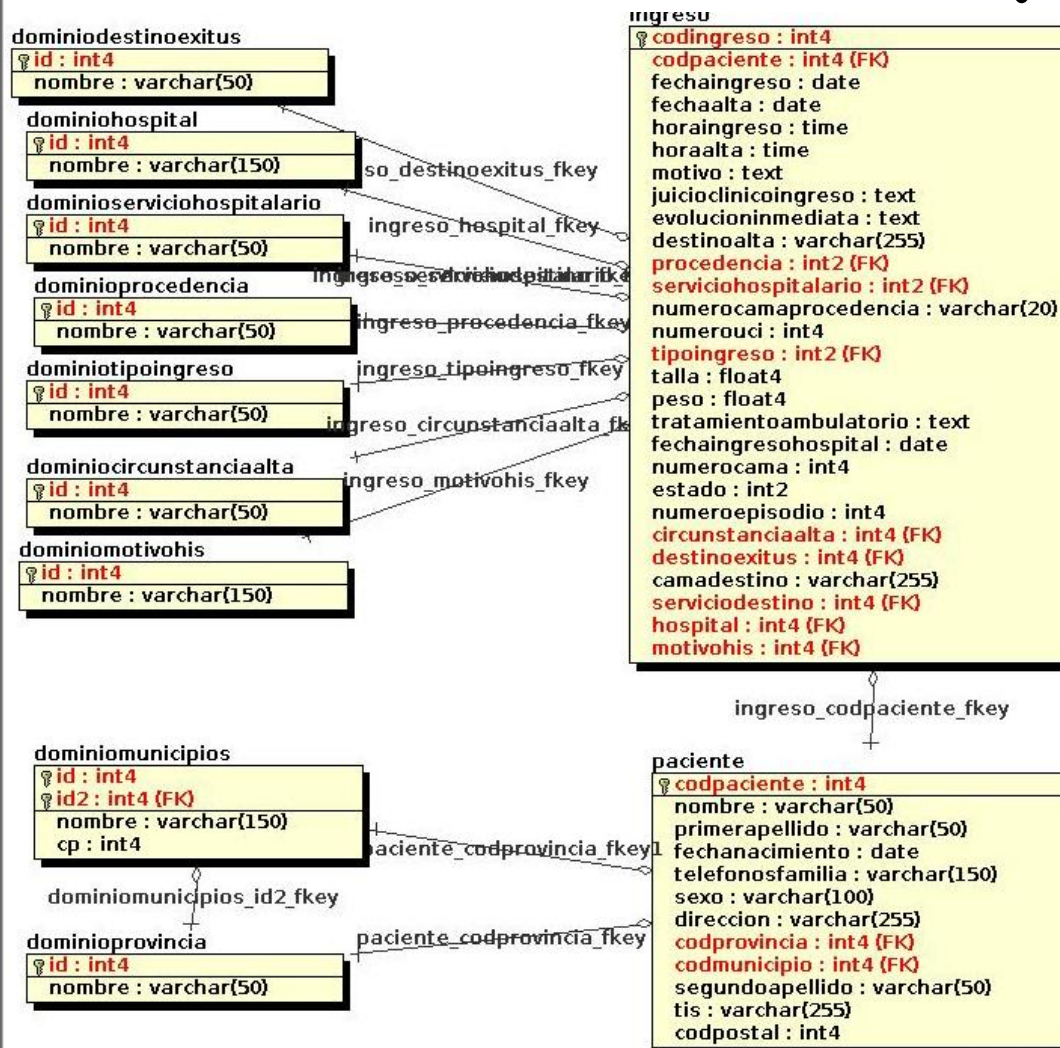
UNIVERSIDAD DE
**MURCIA**

# Pentaho Data Integrator

- Tools

  - Spoon : GUI for transformations and works design

  - Pan : Running transformations

  - Kitchen : Running works (complex transformations)

  - Carte : Cluster Web Server

- Download:

  - Pentaho Data Integrator (v 7.1)

    - http://sourceforge.net/projects/pentaho/files/

  - (if needed Driver Postgresql-9.42, jre8) dir: $PDI/lib/

    - http://jdbc.postgresql.org/download.html

- Tutorial: lets see some simple tasks.

  - http://wiki.pentaho.com/display/EAI/Pentaho+Data+Integration+Steps

- Parallel nature:

    - When a transformation is launched, all its steps are started.

    - During the execution, the steps work simultaneously reading rows from the incoming hops, processing them, and delivering them to the outgoing hops. When there are no more rows left, the execution of the transformation ends.

- Streams:

    - A stream is a set of rows all having the same structure or metadata. (Metadata is important)

    - Split, join, inject, …

- Identify rows when reading (rownum)

- Source schema: pacientes



- Check or create the connection->

  - Menu "Tools"->Wizard
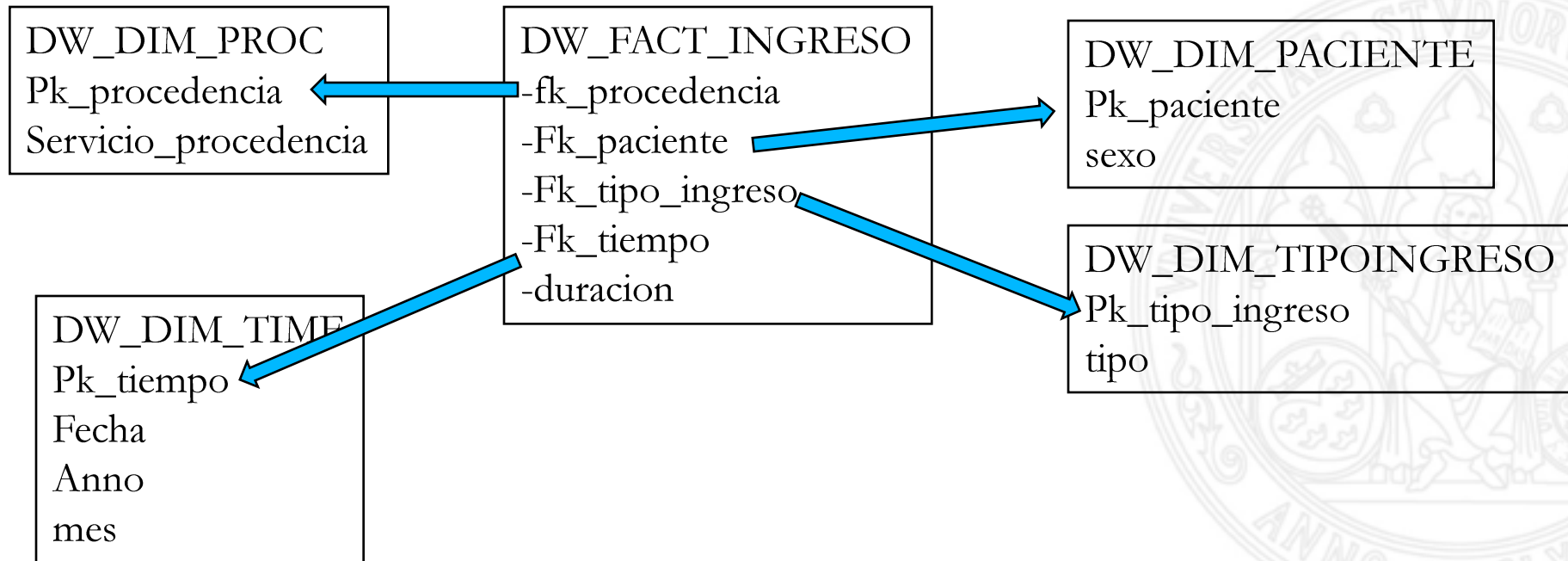
    - "Create database connection"

      - Postgresql
      - Native (JDBC)
      - Server: given:3128
      - bd: tut3
      - User/pwd: given

- **Check: Tools->Database->Explore**

- **Each of you must use your own schema:**

  - **dw_inbdXX**

# Database

- Target schema



**DW_DIM_PROC**
Pk_procedencia
Servicio_procedencia

**DW_FACT_INGRESO**
-fk_procedencia
-Fk_paciente
-Fk_tipo_ingreso
-Fk_tiempo
-duracion

**DW_DIM_PACIENTE**
Pk_paciente
sexo

**DW_DIM_TIPOINGRESO**
Pk_tipo_ingreso
tipo

**DW_DIM_TIME**
Pk_tiempo
Fecha
Anno
mes

# Starting

- Lauch spoon with data_integration/spoon.sh

- Create a database conection

  - Create a first transformation: "File"-> "New"->"Transformation"

  - Menu "Tools" -> "Wizard" -> "Create Connection" and Test it.

  - "View" tab ->Right click on "Database Connection" -> Share: It appear in E/.kettle/shared.xml



**Note that IP changes**

# Load dimension: Procedence

- Create new transformation: Save it as "LoadProcedence"

- Select "Design" tab-> Input->Table Input , and Drag & Drop into the right

  - Select your Connection

  - Define query (we can also browse with "Obtain query").

    - `Select id,nombre from pacientes.dominioprocedencia;`

  - Check with preview

# Load dimension: Procedence

- Select Design → Out → Out table

  - Define **YOUR** target schema "dw_inbd_XX" and table "dw_dim_proc".

  - Note: Batch inserts for efficiency/

  - Also "Insert/Update" step.

# Load dimension: Procedence

- We have to link the stepts:

  - Pressing shift, central button, drag&drop, ….

  - Go to previous step and "WritingProcedences" map the fields with the stream:

  - Choose "Specify database fields"

    - Tab: "Database fields".

  - If needed

    - Use "Transform-> Select values"

    - to rename or metadata

# Run transformation

- **IMPORTANT:** Before running: Verify transformation

- Several running options

- Parameter passing

- Log level

- Safe mode:

  - Extra row checking in runtime

UNIVERSIDAD DE
**MURCIA**



- It can be debugged

  - See "Write to log"

  - Error handling in some steps. Rigth click



11

# Run transformation: Metrics

**Execution Results**

Execution History | Logging | Step Metrics | Performance Graph | Metrics | Preview data

| | Stepname | Copynr | Read | Written | Input | Output | Updated | Rejected | Errors | Active | Time |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | ReadingProcedences | 0 | 0 | 3 | 3 | 0 | 0 | 0 | 0 | Stopped | 0.0s |
| 2 | WritingProcedences | 0 | 3 | 0 | 0 | 0 | 0 | 0 | 1 | Stopped | 0.0s |

**Execution Results**

Execution History | Logging | Step Metrics | Performance Graph | Metrics | Preview data

```
                5        10        15        20        25        30        35

Execute a transformation - cargarProcedencia : 39ms

Initialize a transformation - cargarProcedencia : 12ms

    Connect to database - tut3 : 10ms

    Initialize a step - ReadingProcedences : 10ms

    Connect to database - tut3 : 11ms

    Initialize a step - WritingProcedences : 11ms

                              Open SQL query - tut3 : 11ms

                              Execute SQL statement - tut3 : 2ms

                              Create SQL statement - tut3 : 0ms

                              Execute a step - ReadingProcedences : 17ms

                              Execute a step - WritingProcedences : 24ms

                                  Get row metadata - tut3 : 9ms

                                            Get DB metadata - tut3 : 0ms
```
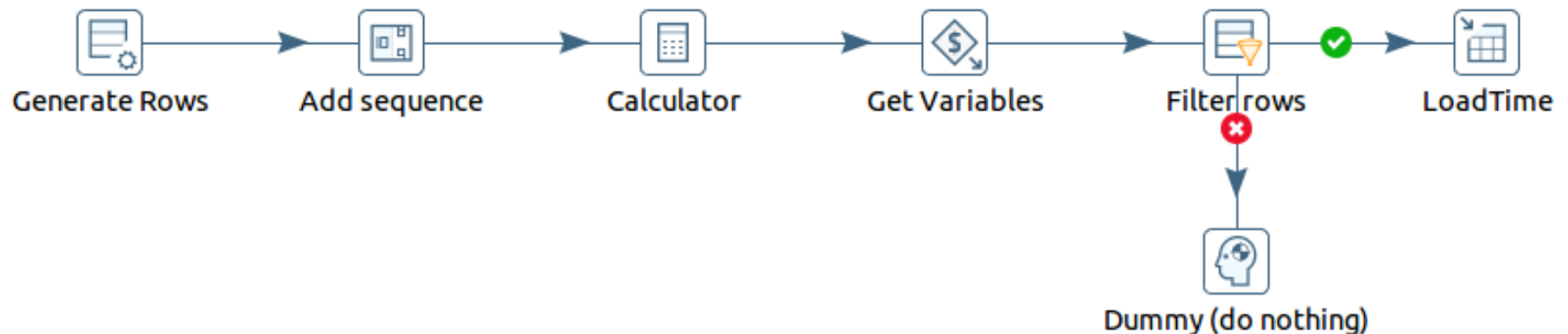
12

# Load dimensions

- We do the same for

  - Admission type
    - Source table: pacientes.dominiotipoingreso
    - Target table: dw_YOUR_SCHEMA.dw_dim_tipoingreso
    - Map both fields
    - *Save it as "LoadAdmissionType"*

  - Patient:
    - Source table: pacientes.paciente
      - Source fields: "codpaciente", "sexo"
    - Target table: dw_YOUR_SCHEMA.dw_dim_patient
    - Map both fields
    - *Save it as "LoadPatient"*

# Load time dimension

- Objective (step):

  - Generate 1000 dates (generate rows)

  - Create a sequence to be added (sequence)

  - Add to every date the i-th number of the sequence (calculator)

  - Extract year and month (calculator)

  - Create subrogate key (calculator) -> the same sequence

  - Read parameter (read variable)(define the variable in the transformation)

  - Filter rows older than the parameter date

  - Store in the "time" table using the consecutive pk

# Load time dimension

- We enter the following steps

  1. Input – Generate rows

  2. Transform – Add Sequence

  3. Transform – Calculator

  4. Job – Get variables

  5. Flow – Filter rows

  6. Flow – Dummy transformation

  7. Out – Table writer

# Load time dimension

- Generate rows

  1. Input – Generate rows

  2. Set format and initial value

  3. Limit the number of rows: 800

UNIVERSIDAD DE
**MURCIA**

- Add a day

1. Transform – Add Sequence

2. Set a name for the value

3. Set inital, increment and max value

- Transform – Calculator

1. Calculate new date

2. Extract year and month

3. We also create the subrogate key (comment from lecturer)

4. Note:

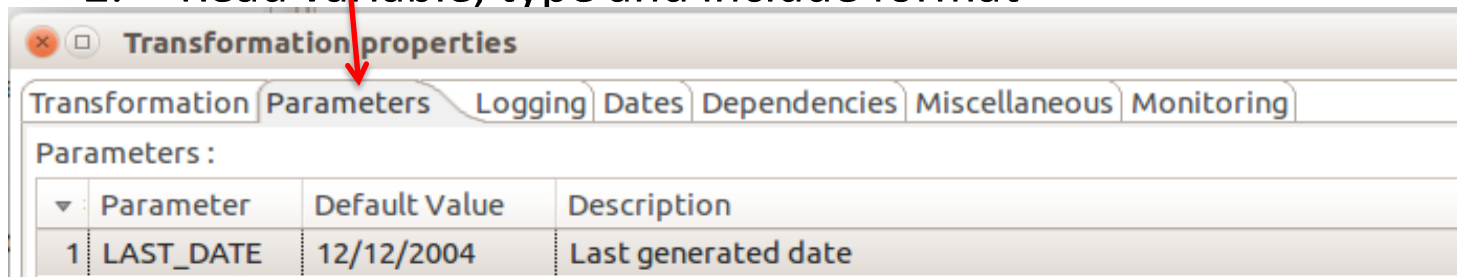   1. Important: Date type

   2. Remove?



**Calculator**

Step name | CalculateAndExtractFields

Fields:

| | New field | Calculation | Field A | Field B | Field C | Value type | Length | Precision | Remove | Conve |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | calculated_date | Date A + B Day | myDate | added_day | | Date | | | N | |
| 2 | year | Year of date A | calculated_date | | | Integer | | | N | |
| 3 | month | Month of date | calculated_date | | | Integer | | | N | |
| 4 | datekey | Create a copy ( | added_day | | | Integer | | | N | |

Help | OK | Cancel

# Load time dimension

- Job – Read variable

  1. (background of transformation) Right Click->Properties: Define in "Parameters" tab the variable and value in the transformation

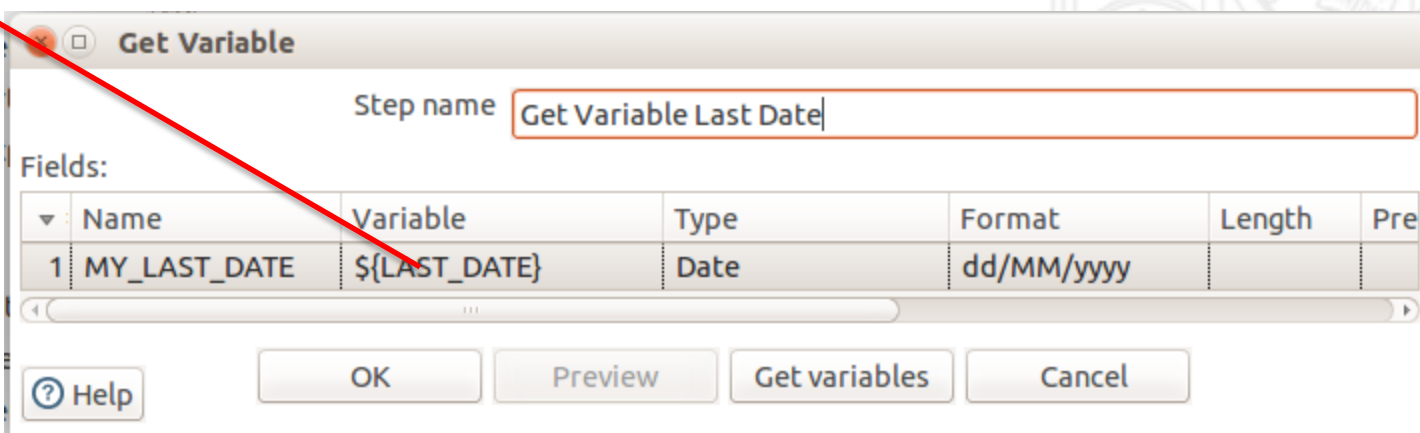  2. Read variable, type and include format

**Transformation properties**

| Transformation | Parameters | Logging | Dates | Dependencies | Miscellaneous | Monitoring |

Parameters :

| | Parameter | Default Value | Description |
|---|---|---|---|
| 1 | LAST_DATE | 12/12/2004 | Last generated date |

**Get Variable**

Step name | Get Variable Last Date

Fields:

| | Name | Variable | Type | Format | Length | Pre |
|---|---|---|---|---|---|---|
| 1 | MY_LAST_DATE | ${LAST_DATE} | Date | dd/MM/yyyy | | |

| ⑦ Help | OK | Preview | Get variables | Cancel |

19

# Load time dimension

- Filter rows

1. Flow – Filtrer rows

2. Flow – Dummy transformation: "Discard dates"

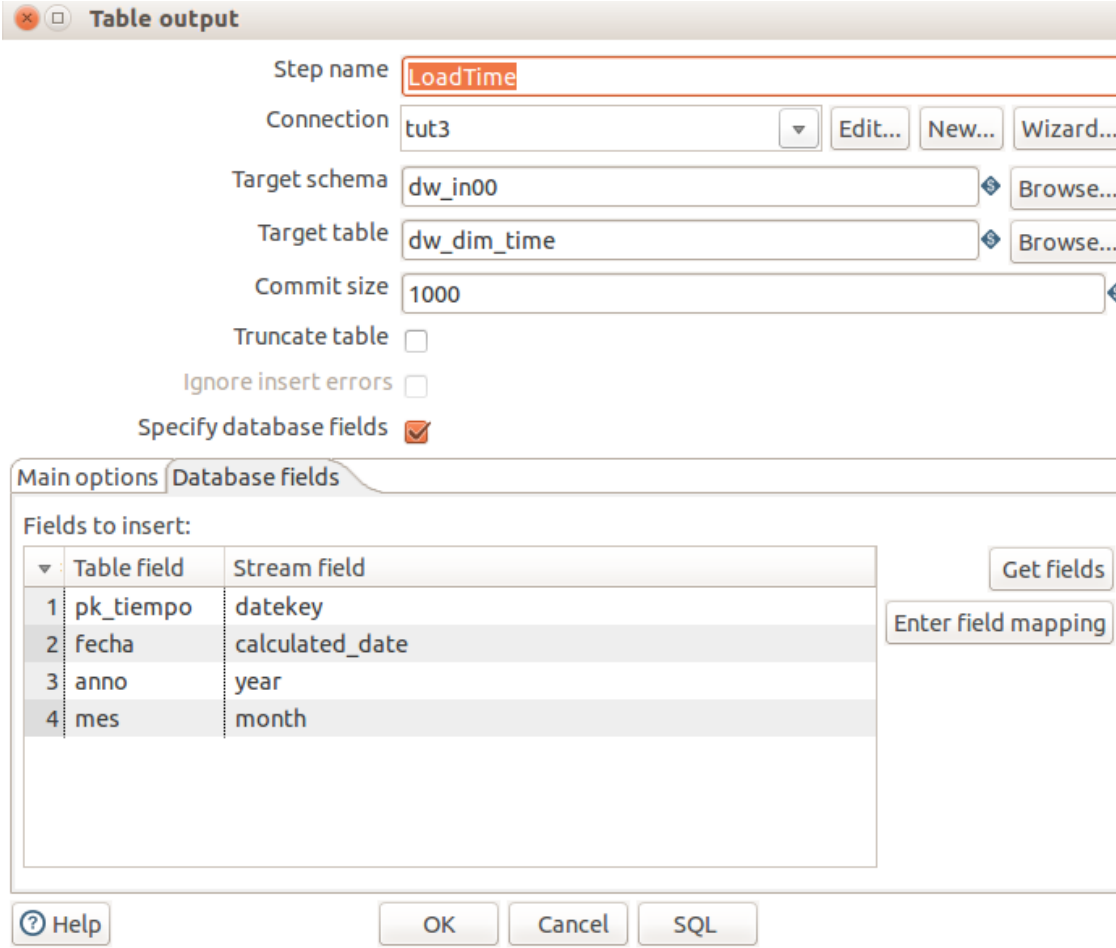3. Sent to dummy transformation rows after last date

UNIVERSIDAD DE
**MURCIA**

- Out – Table writter

  1. Mapping database fields



21

# Load time dimension

- ## Run job with parameters

**Run Options**

**Environment Type**

- ⦿ Local — The transformation will run on the machine you are using.
- ○ Server
- ○ Clustered

**Options**

- ☑ Clear log before running   Log level: Basic
- ☐ Enable safe mode
- ☑ Gather performance metrics

**Parameters** | Variables

| Parameter | Default value | Value | Description |
|-----------|---------------|-------|-------------|
| LAST_DATE | 12/12/2004 |  | Last generated date |

☑ Always show dialog on run

⑦ Help

**Execution Results**

⏱ Execution History | 📋 Logging | 1≡ Step Metrics | ⤢ Performance Graph | ⧉ Metrics | ◉ Preview data

◉

| | Stepname | Copynr | Read | Written | Input | Output | Updated | Rejected | Errors |
|---|----------|--------|------|---------|-------|--------|---------|----------|--------|
| 1 | Generate Rows | 0 | 0 | 800 | 0 | 0 | 0 | 0 | 0 |
| 2 | Add day | 0 | 800 | 800 | 0 | 0 | 0 | 0 | 0 |
| 3 | CalculateAndExtractFields | 0 | 800 | 800 | 0 | 0 | 0 | 0 | 0 |
| 4 | Get Variable Last Date | 0 | 800 | 800 | 0 | 0 | 0 | 0 | 0 |
| 5 | Filter valid dates | 0 | 800 | 800 | 0 | 0 | 0 | 0 | 0 |
| 6 | Dummy (do nothing) | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 7 | LoadTime | 0 | 800 | 800 | 0 | 800 | 0 | 0 | 0 |

- Lets do the following

  1. Input -> Input table: Select some admissions

  2. Lookoup -> Database lookup: Foreign key in the time table

  3. Ouput -> Ouput table: Store the facts

- Note: this is not the normal procedure, just an example, we are not creating subrogate keys except for time.

ReadAdmissions      Database lookup      WriteAdmissions

# Load facts

- In fact table load:

1. codpaciente

2. procedencia

3. tipoingreso

4. fechaingreso,

5. Meassure: duration

Conditions:

1. Discharged

   - fechaalta is not null

2. Correct dates

   - discharge >= admission

3. Correct tipoingreso

   - tipoingreso is not null

# Load facts

- Look up the primary key in time dimension matching the admission date

    1. Search – Database lookup

    2. Select lookup fields

    3. Select extracted field:

        1. Dimension pk

    - Remaining fields: optional

    - Careful with type

# Load facts

- Store output in table
  - Map the fields

# Create job

- ## Put all transformations in a job

  - ### New File -> Job

  - ### General -> Start

  - ### General ->Transformation

  - ### General -> Success

- Condition for links

- We can add a new transformation for creating or rebuilding the database

- Include also an error node to abort the job

# Interesting Steps

- Input/output: CSV, XML, datagrid, Property files

- Transform: Select values

- Scripting: SQL Scripting, JavaScript

- Lookup: Database join

- Mapping: to create "subtransformations"

- Job: Get/set files/variables from results

- Filemanagement, filetransfer, mail (for jobs)

- Upload to the task in Aula Virtual a compressed file with the following:

  - Transformation definitions (.ktr files)

    - Load procedence
    - Load admission type
    - Load patient
    - Load time
    - Load Fact tables
    - Create and drop schema

  - Job definition (.kjb file)

    - Load all