# Data Analysis and Integration

ETL tools and SQL queries

# Pentahoo Data Integration

- Transformations work in "streaming" mode
  - the first row could reach the output before the second row is read from input
  - there are exceptions, e.g. when some aggregation on multiple rows needs to be performed
- Transformations can be saved and executed many times
- Transformations can become quite complex

# SQL query

- number of employees in each department

curr_dept_emp

```
+---------+---------+
| emp_no  | dept_no |
+---------+---------+
|   10721 | d009    |
|   11260 | d009    |
|   11371 | d005    |
|   11693 | d005    |
|   13816 | d005    |
|   14007 | d002    |
|   14083 | d004    |
|   14791 | d005    |
|   17698 | d005    |
|   17739 | d005    |
|     ... |   ...   |
+---------+---------+
```

departments

```
+---------+--------------------+
| dept_no | dept_name          |
+---------+--------------------+
| d009    | Customer Service   |
| d005    | Development        |
| d002    | Finance            |
| d003    | Human Resources    |
| d001    | Marketing          |
| d004    | Production         |
| d006    | Quality Management |
| d008    | Research           |
| d007    | Sales              |
+---------+--------------------+
```

# SQL query

- number of employees in each department

**select** *
**from** curr_dept_emp **as** a, departments **as** b
**where** a.dept_no = b.dept_no;

```
+--------+---------+---------+------------------+
| emp_no | dept_no | dept_no | dept_name        |
+--------+---------+---------+------------------+
|  10721 | d009    | d009    | Customer Service |
|  11260 | d009    | d009    | Customer Service |
|  11371 | d005    | d005    | Development      |
|  11693 | d005    | d005    | Development      |
|  13816 | d005    | d005    | Development      |
|  14007 | d002    | d002    | Finance          |
|  14083 | d004    | d004    | Production       |
|  14791 | d005    | d005    | Development      |
|  17698 | d005    | d005    | Development      |
|  17739 | d005    | d005    | Development      |
+--------+---------+---------+------------------+
10 rows in set (0.00 sec)
```

# SQL query

- number of employees in each department

**select** b.dept_no, b.dept_name, **count**(emp_no)
**from** curr_dept_emp **as** a, departments **as** b
**where** a.dept_no = b.dept_no
**group by** b.dept_no, b.dept_name;

```
+---------+--------------------+---------------+
| dept_no | dept_name          | count(emp_no) |
+---------+--------------------+---------------+
| d001    | Marketing          |            15 |
| d002    | Finance            |            18 |
| d003    | Human Resources    |            10 |
| d004    | Production         |            44 |
| d005    | Development        |            62 |
| d006    | Quality Management |            18 |
| d007    | Sales              |            42 |
| d008    | Research           |            14 |
| d009    | Customer Service   |            29 |
+---------+--------------------+---------------+
9 rows in set (0.00 sec)
```

# SQL query

- number of employees in each department

**select** b.dept_no, b.dept_name, **count**(emp_no) **as** count_emp_no
**from** curr_dept_emp **as** a, departments **as** b
**where** a.dept_no = b.dept_no
**group by** b.dept_no, b.dept_name
**having** count_emp_no >= 40;

```
+---------+-------------+--------------+
| dept_no | dept_name   | count_emp_no |
+---------+-------------+--------------+
| d004    | Production  |           44 |
| d005    | Development |           62 |
| d007    | Sales       |           42 |
+---------+-------------+--------------+
3 rows in set (0.00 sec)
```

# SQL query

- number of employees in each department

**select** b.dept_no, b.dept_name, **count**(emp_no) **as** count_emp_no
**from** curr_dept_emp **as** a, departments **as** b
**where** a.dept_no = b.dept_no
**group by** b.dept_no, b.dept_name
**having** count_emp_no >= 40
**order by** count_emp_no **desc**;

```
+---------+-------------+--------------+
| dept_no | dept_name   | count_emp_no |
+---------+-------------+--------------+
| d005    | Development |           62 |
| d004    | Production  |           44 |
| d007    | Sales       |           42 |
+---------+-------------+--------------+
3 rows in set (0.00 sec)
```

# SQL query

- doing the same query with an ETL tool

    **select** b.dept_no, b.dept_name, **count**(emp_no) **as** count_emp_no
    **from** curr_dept_emp **as** a, departments **as** b
    **where** a.dept_no = b.dept_no
    **group by** b.dept_no, b.dept_name
    **having** count_emp_no >= 40
    **order by** count_emp_no **desc**;

# –how to do it?

# Table input

**select** a.dept_no, b.dept_name, **count**(a.emp_no) **as** count_emp_no
**from** curr_dept_emp **as** a, departments **as** b
**where** a.dept_no = b.dept_no
**group by** a.dept_no, b.dept_name
**having** count_emp_no >= 40
**order by** count_emp_no **desc**;

# Table input

**select** a.dept_no, b.dept_name, **count**(a.emp_no) **as** count_emp_no
**from** curr_dept_emp **as** a, departments **as** b
**where** a.dept_no = b.dept_no
**group by** a.dept_no, b.dept_name
**having** count_emp_no >= 40
**order by** count_emp_no **desc**;



Table input 2

```
Table input

Step name    Table input 2
Connection   employees

SQL

SELECT
  dept_no
, dept_name
FROM departments
```

# Select values (for renaming)

**select** a.dept_no, b.dept_name, **count**(a.emp_no) **as** count_emp_no
**from** curr_dept_emp **as** a, departments **as** b
**where** a.dept_no = b.dept_no
**group by** a.dept_no, b.dept_name
**having** count_emp_no >= 40
**order by** count_emp_no **desc**;

# Select values (for renaming)

**select** a.dept_no, b.dept_name, **count**(a.emp_no) **as** count_emp_no
**from** curr_dept_emp **as** a, departments **as** b
**where** a.dept_no = b.dept_no
**group by** a.dept_no, b.dept_name
**having** count_emp_no >= 40
**order by** count_emp_no **desc**;

# Join rows

**select** a.dept_no, b.dept_name, **count**(a.emp_no) **as** count_emp_no
**from** curr_dept_emp **as** a, departments **as** b
**where** a.dept_no = b.dept_no
**group by** a.dept_no, b.dept_name
**having** count_emp_no >= 40
**order by** count_emp_no **desc**;

# Sort rows

**select** a.dept_no, b.dept_name, **count**(a.emp_no) **as** count_emp_no
**from** curr_dept_emp **as** a, departments **as** b
**where** a.dept_no = b.dept_no
**group by** a.dept_no, b.dept_name
**having** count_emp_no >= 40
**order by** count_emp_no **desc**;

# Group by

**select** a.dept_no, b.dept_name, **count**(a.emp_no) **as** count_emp_no
**from** curr_dept_emp **as** a, departments **as** b
**where** a.dept_no = b.dept_no
**group by** a.dept_no, b.dept_name
**having** count_emp_no >= 40
**order by** count_emp_no **desc**;

# Group by

select a.dept_no, b.dept_name, **count**(a.emp_no) **as** count_emp_no
**from** curr_dept_emp **as** a, departments **as** b
**where** a.dept_no = b.dept_no
**group by** a.dept_no, b.dept_name
**having** count_emp_no >= 40
**order by** count_emp_no **desc**;



🗗ₓ Group By

🗙 Warning!                                                    ✕

⚠ The group by function needs the input to be sorted on the specified keys.
If you don't sort the input, the results may not be correct

☐ Please, don't show this warning anymore.

[ I understand ]

The fields that make up the group:

| # | Group field |
|---|---|
| 1 | dept_no_1 |
| 2 | dept_name |

Aggregates :

| # | Name | Subject | Type |
|---|---|---|---|
| 1 | count_emp_no | emp_no | Number of Values (N) |

# Filter rows

**select** a.dept_no, b.dept_name, **count**(a.emp_no) **as** count_emp_no
**from** curr_dept_emp **as** a, departments **as** b
**where** a.dept_no = b.dept_no
**group by** a.dept_no, b.dept_name
**having** count_emp_no >= 40
**order by** count_emp_no **desc**;

# Filter rows

**select** a.dept_no, b.dept_name, **count**(a.emp_no) **as** count_emp_no
**from** curr_dept_emp **as** a, departments **as** b
**where** a.dept_no = b.dept_no
**group by** a.dept_no, b.dept_name
**having** count_emp_no >= 40
**order by** count_emp_no **desc**;

# Preview

# Preview

# Preview

# Preview

# Preview

# Preview

# Preview

# Preview

# Preview

# Alternatives

- A completely different solution



Table input

# Alternatives

- A completely different solution

# Alternatives

- Other possible intermediate solutions
  - column renaming in SQL



  - column renaming and table join in SQL



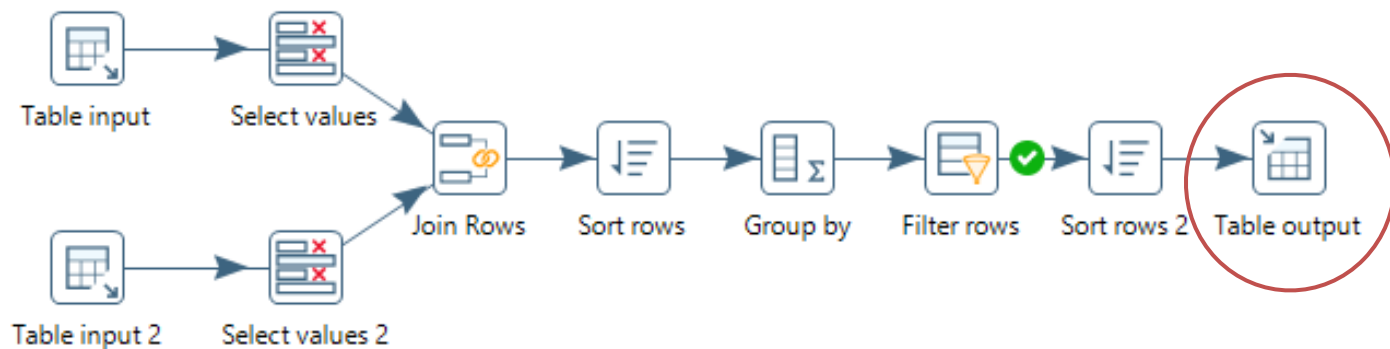  - column renaming, table join, and group by in SQL

# Alternatives

- As a general rule

  - **offload as much as possible to database system**
    - filtering rows, selecting columns, and (especially) table joins

  - **unless the required data comes from multiple data sources**

  - **or when processing is difficult to implement in SQL**
    - e.g. duplicate detection via string matching

# ETL tools

- Then why use an ETL tool?
    - data comes from different databases/systems
    - data sources other than databases (e.g. text files)
    - complex data merging and transformations
    - approximate matching, duplicate detection, data cleaning
    - materialization to different outputs (databases, files, etc.)

# Materialization

- Materialization to database table

# Materialization

- Materialization to database table

# Materialization

- Materialization to database table

# Materialization

- Running the transformation

# Materialization

- Running the transformation

**Execution Results**

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

| # | Stepname | Copynr | Read | Written | Input | Output | Updated | Rejected | Errors | Active | Time | Speed (r/s) |
|---|----------|--------|------|---------|-------|--------|---------|----------|--------|--------|------|-------------|
| 1 | Table input | 0 | 0 | 252 | 252 | 0 | 0 | 0 | 0 | Finished | 0.0s | 19,385 |
| 2 | Table input 2 | 0 | 0 | 9 | 9 | 0 | 0 | 0 | 0 | Finished | 0.0s | 1,800 |
| 3 | Select values | 0 | 252 | 252 | 0 | 0 | 0 | 0 | 0 | Finished | 0.0s | 14,824 |
| 4 | Select values 2 | 0 | 9 | 9 | 0 | 0 | 0 | 0 | 0 | Finished | 0.0s | 818 |
| 5 | Join Rows | 0 | 261 | 252 | 0 | 0 | 0 | 0 | 0 | Finished | 0.3s | 821 |
| 6 | Sort rows | 0 | 252 | 252 | 0 | 0 | 0 | 0 | 0 | Finished | 0.3s | 775 |
| 7 | Group by | 0 | 252 | 9 | 0 | 0 | 0 | 0 | 0 | Finished | 0.3s | 766 |
| 8 | Filter rows | 0 | 9 | 3 | 0 | 0 | 0 | 0 | 0 | Finished | 0.3s | 27 |
| 9 | Sort rows 2 | 0 | 3 | 3 | 0 | 0 | 0 | 0 | 0 | Finished | 0.3s | 9 |
| 10 | Table output | 0 | 3 | 3 | 0 | 3 | 0 | 0 | 0 | Finished | 0.3s | 9 |

# Materialization

- Materialization to database table

**select** * **from** num_emp_dept;

```
+-----------+-------------+--------------+
| dept_no_1 | dept_name   | count_emp_no |
+-----------+-------------+--------------+
| d005      | Development |           62 |
| d004      | Production  |           44 |
| d007      | Sales       |           42 |
+-----------+-------------+--------------+
3 rows in set (0.00 sec)
```