



PostgreSQL

PGDay.IT 2011

*Monash University Prato Centre
Venerdì 25 Novembre 2011*

ETL con Kettle

Giulio Calacoci

Italian PostgreSQL Users Group

www.itpug.org

www.postgresql.org



Chi sono?

- 2ndQuadrant Italia:
 - Business Intelligence
 - Software Development:
 - Java
 - PHP



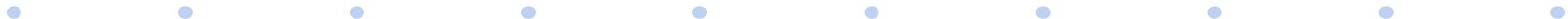
Sommario

- ETL
 - Concetti di base
 - Esempi di applicazione
- Kettle
 - Componenti
 - Casi d'uso reale



Principali Utilizzatori

- Developer:
 - Migrazioni di Applicazioni
 - Migrazioni di Database
 - Importazione / Esportazione dati
 - QA
- Business Intelligence Expert:
 - OLTP
 - Fonti di dati ausiliarie
 - Webservice
 - QA

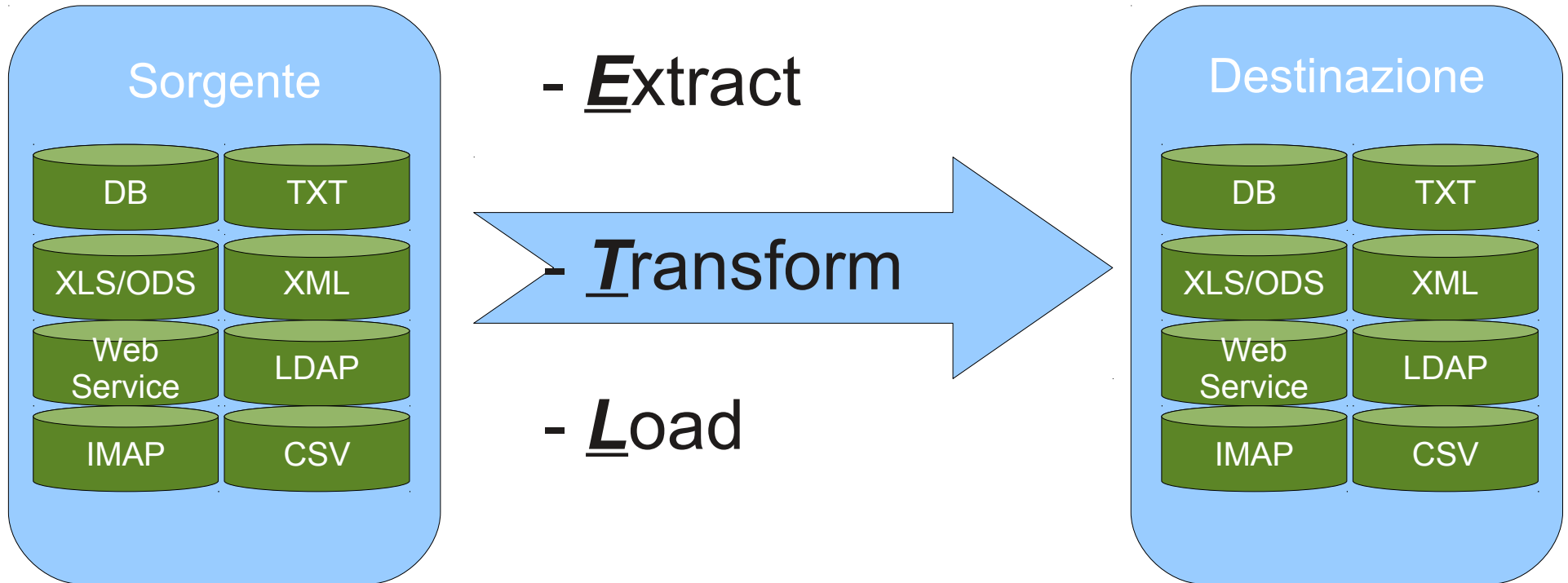


ETL

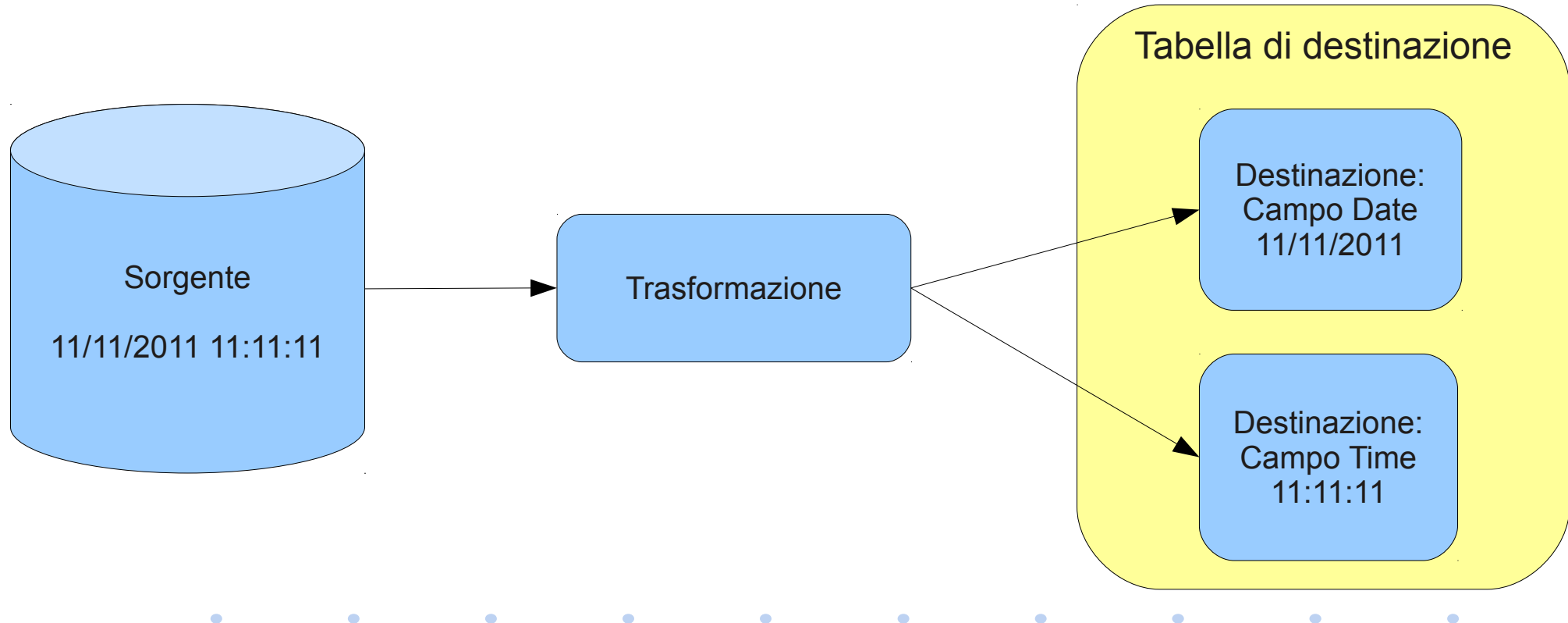
Chi non ha mai scritto un programma del genere?

```
<?php
function csv_file_to_mysql_table($source_file, $target_table, $max_line_length=10000) {
    if (($handle = fopen("$source_file", "r")) !== FALSE) {
        $columns = fgetcsv($handle, $max_line_length, ",");
        foreach ($columns as &$column) {
            $column = str_replace(".", "", $column);
        }
        $insert_query_prefix = "INSERT INTO $target_table (".join(",",$columns).")\nVALUES";
        while (($data = fgetcsv($handle, $max_line_length, ",")) !== FALSE) {
            while (count($data)<count($columns))
                array_push($data, NULL);
            $query = "$insert_query_prefix (".join(",",quote_all_array($data)).").";";
            mysql_query($query);
        }
        fclose($handle);
    }
}
```

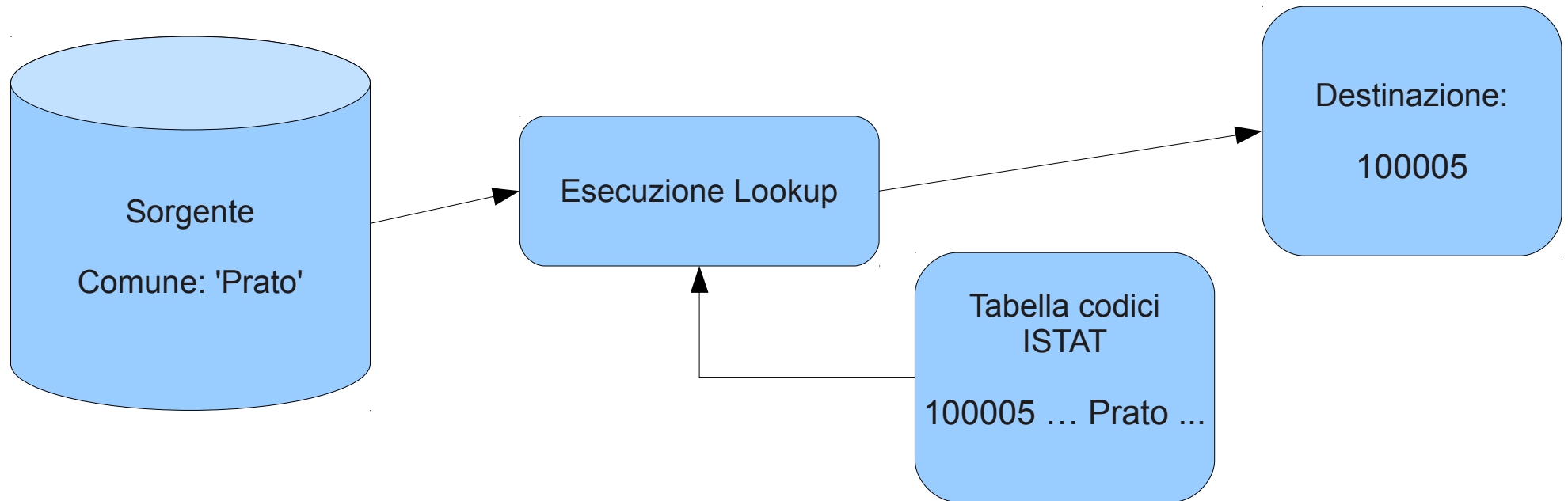
ETL



Esempio: Manipolazione dati



Esempio: Importazione con Lookup



Kettle

- Open Source
 - GNU GPL
 - Pentaho Data Integration
- Scritto in Java
- OLTP (PostgreSQL)
- Warehouse (PostgreSQL o Greenplum)

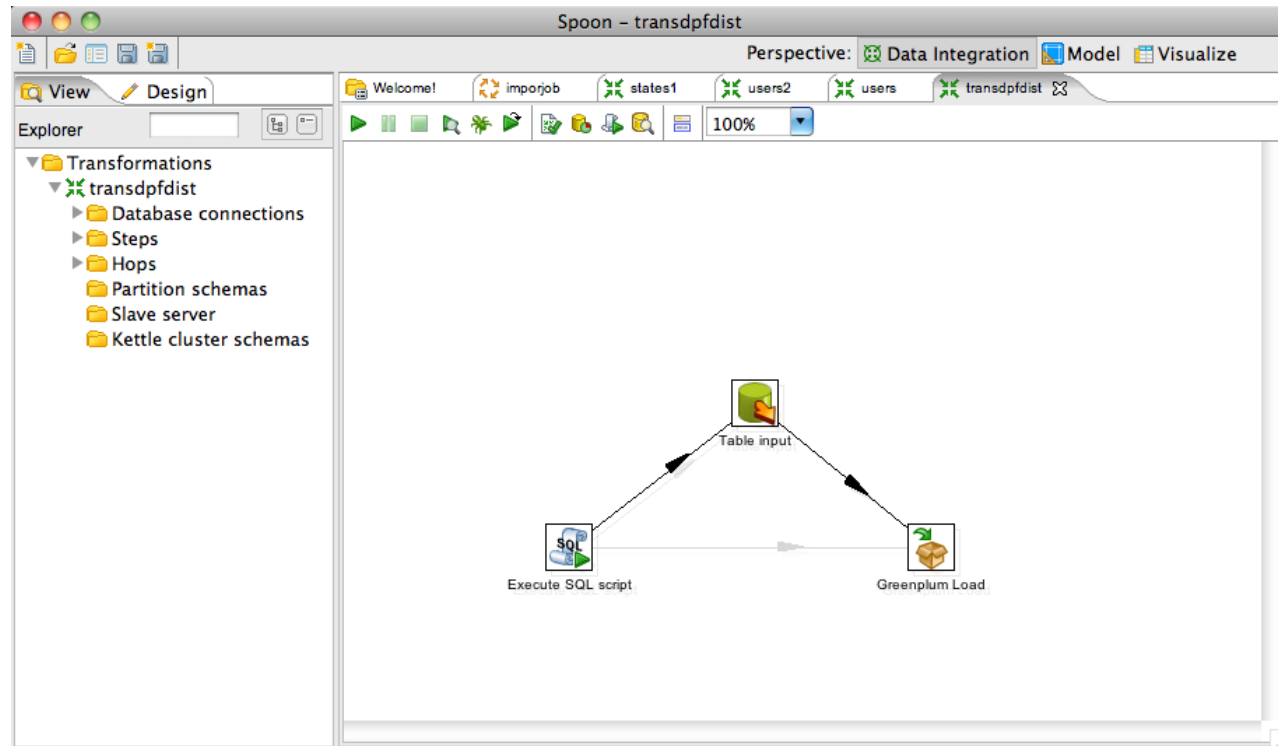
Composto da :

- SPOON
- PAN
- KITCHEN



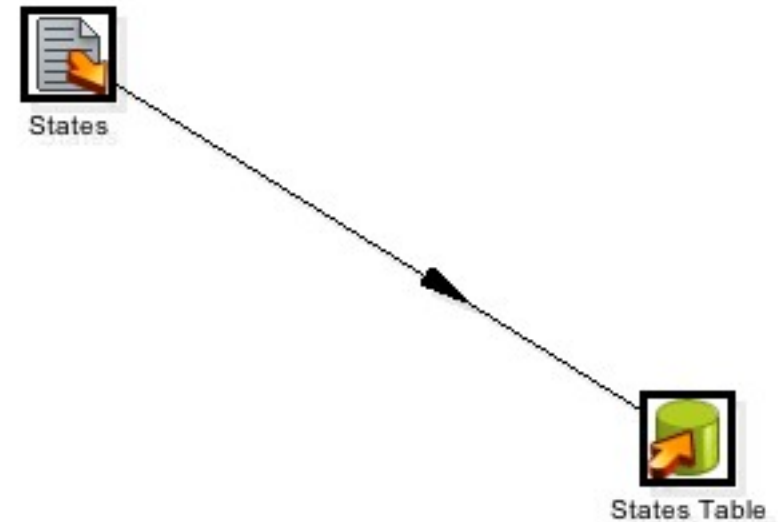
SPOON

- Interfaccia grafica
- Modellazione
- Flusso dati
- Niente codice generato
- Job
- Trasformazioni



Trasformazione

- Sottoinsieme delle operazioni che compongono un job
- Possibilità di esecuzione in parallelo
- Riutilizzabile
- Insieme di componenti specifici per l'esecuzione delle azioni



Componenti

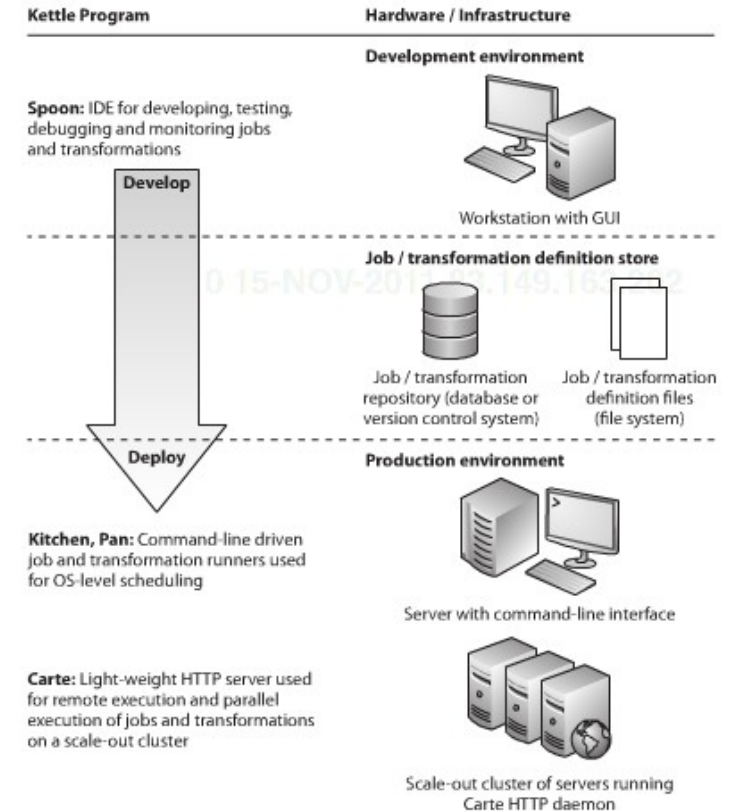
- Organizzati in cartelle
- Rappresentano lo “step” di una trasformazione
- Eseguono operazioni specifiche
- Facilmente configurabili

▼ Bulk loading

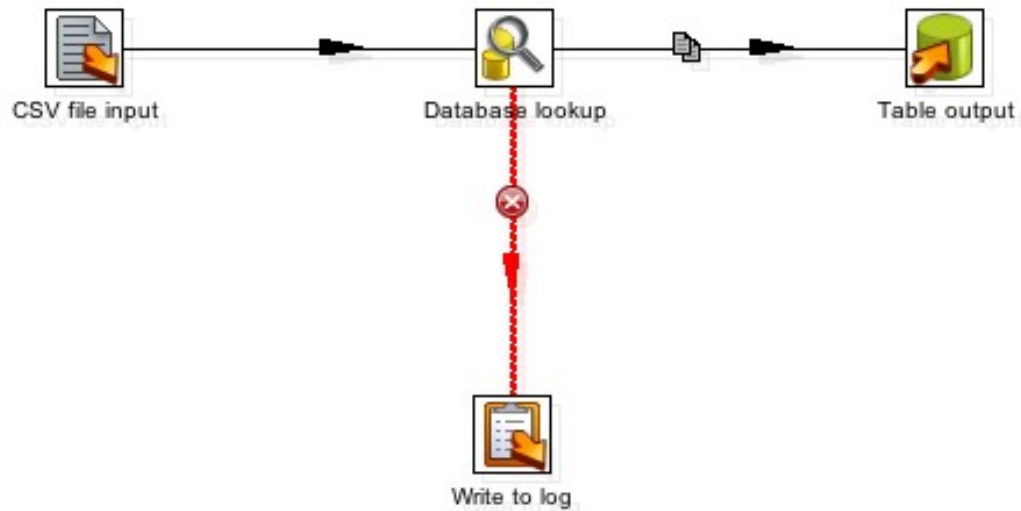
- 🗄️ ElasticSearch Bulk Insert
- 🌱 Greenplum Bulk Loader
- 🌱 Greenplum Load
- 📊 Infobright Loader
- 🌱 Ingres VectorWise Bulk Loader
- 🌱 LucidDB Streaming Loader
- 🌱 MonetDB Bulk Loader
- 🌱 MySQL Bulk Loader
- 🌱 Oracle Bulk Loader
- 🌱 PostgreSQL Bulk Loader
- 🌱 Teradata Fastload Bulk Loader

Automazione

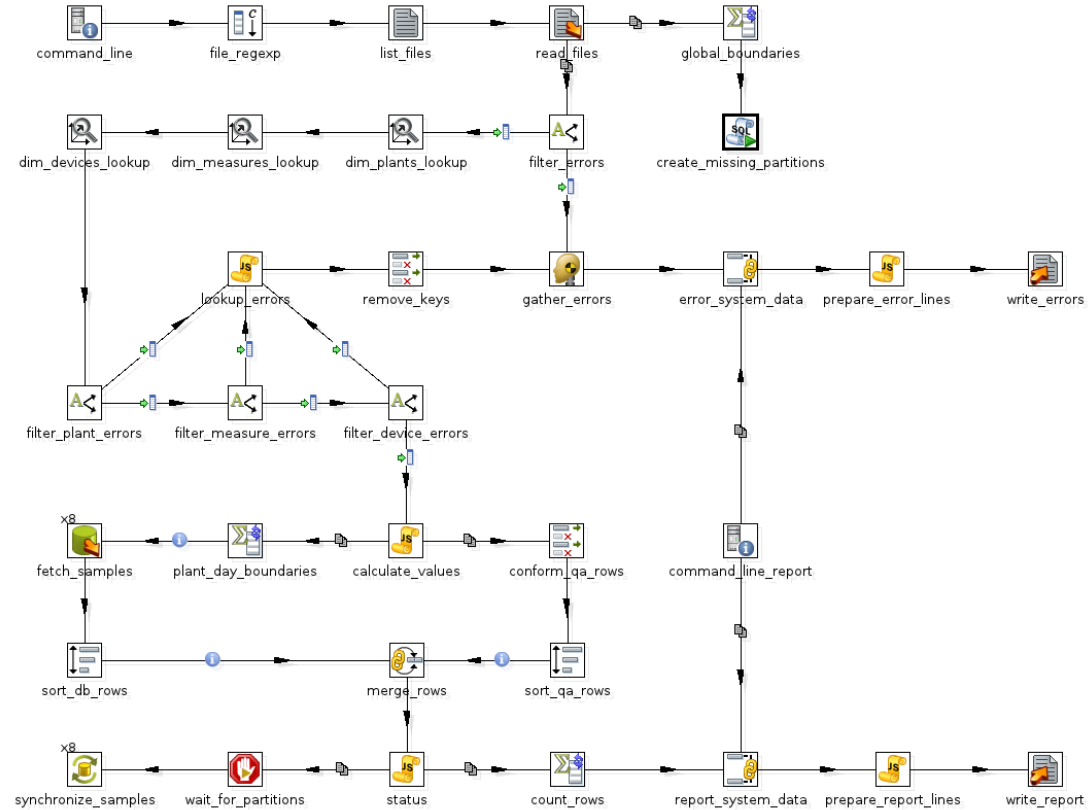
- Fase immediatamente successiva al design
- Eseguibile su server tramite componenti della suite dedicati: Kitchen/Pan
- Esecuzione dei job via cron



Esempi di casi reali: Lookup

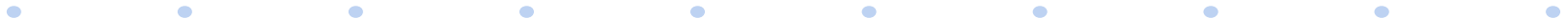


Esempio di caso reale: QA Complesso



Conclusioni

- Figura professionale BI Expert
- Estensibile
- Semplificazione operazioni complesse
- Portabile
- Scalabile
- Cache interna



Domande?

- E-Mail: giulio.calacoci@2ndquadrant.it
- URL: www.2ndquadrant.it
- Blog : blog.2ndquadrant.com



Licenza Creative Commons

Attribuzione

Non commerciale

Condividi allo stesso modo

2.5 Italia

<http://creativecommons.org/licenses/by-nc-sa/2.5/it/>

Copyright 2011 2ndQuadrant Italia - www.2ndQuadrant.it

