

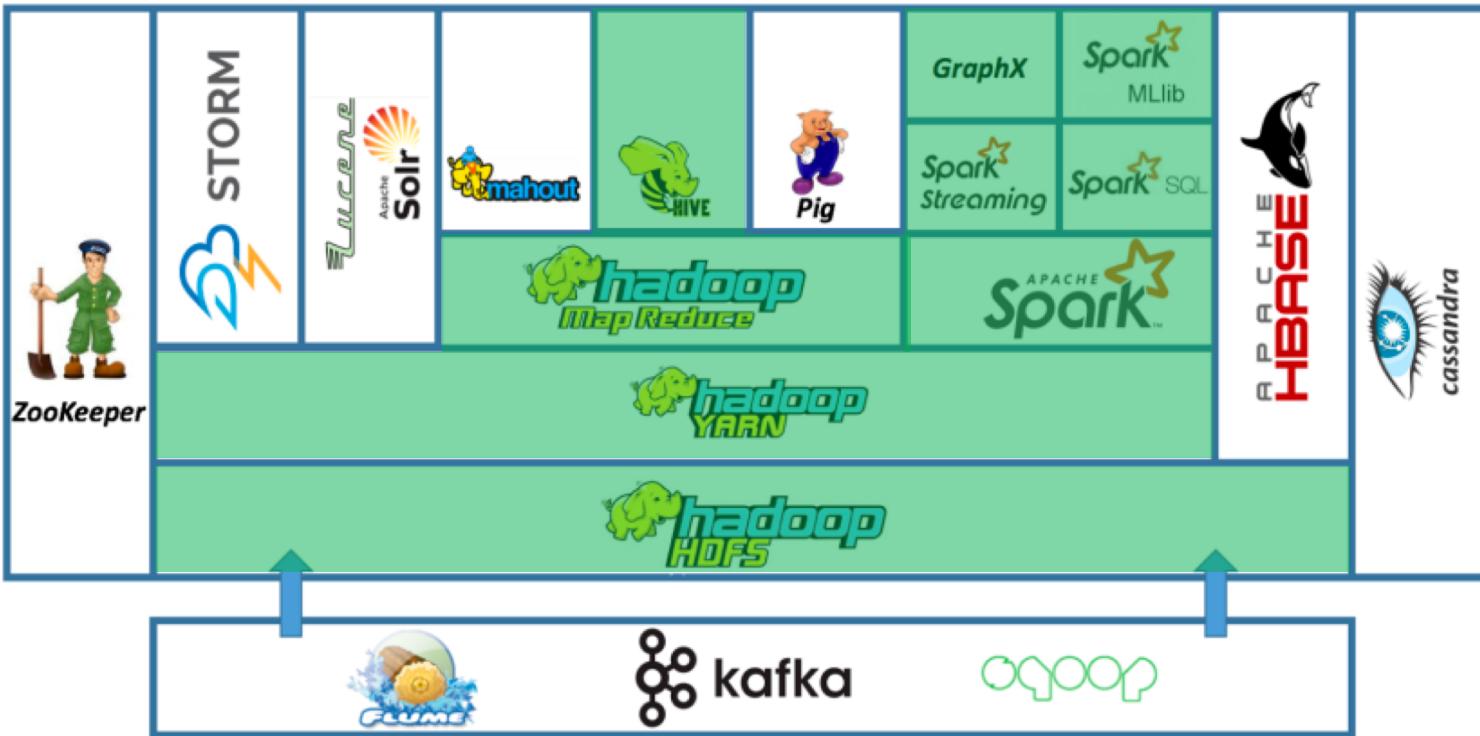
# HandsOn – Pentaho Data Integration

A quick Introduction to ETL Workflow with Pentaho Data Integration



[www.marcel-mittelstaedt.com](http://www.marcel-mittelstaedt.com)

# The Hadoop Ecosystem



Today's  
(exercise) focus



# Exercises Preparation I

Install and Setup Pentaho Data Integration



[www.marcel-mittelstaedt.com](http://www.marcel-mittelstaedt.com)

# Download And Install PDI

## 1. Download Pentaho Data Integration (8.1.0)

```
https://sourceforge.net/projects/pentaho/
```

## 2. Extract and Move:

```
unzip pdi-ce-8.1.0.0-365.zip
```

```
mv data-integration /opt/pdi
```



# Start PDI (Spoon)

## 3. Start Pentaho Data Integration

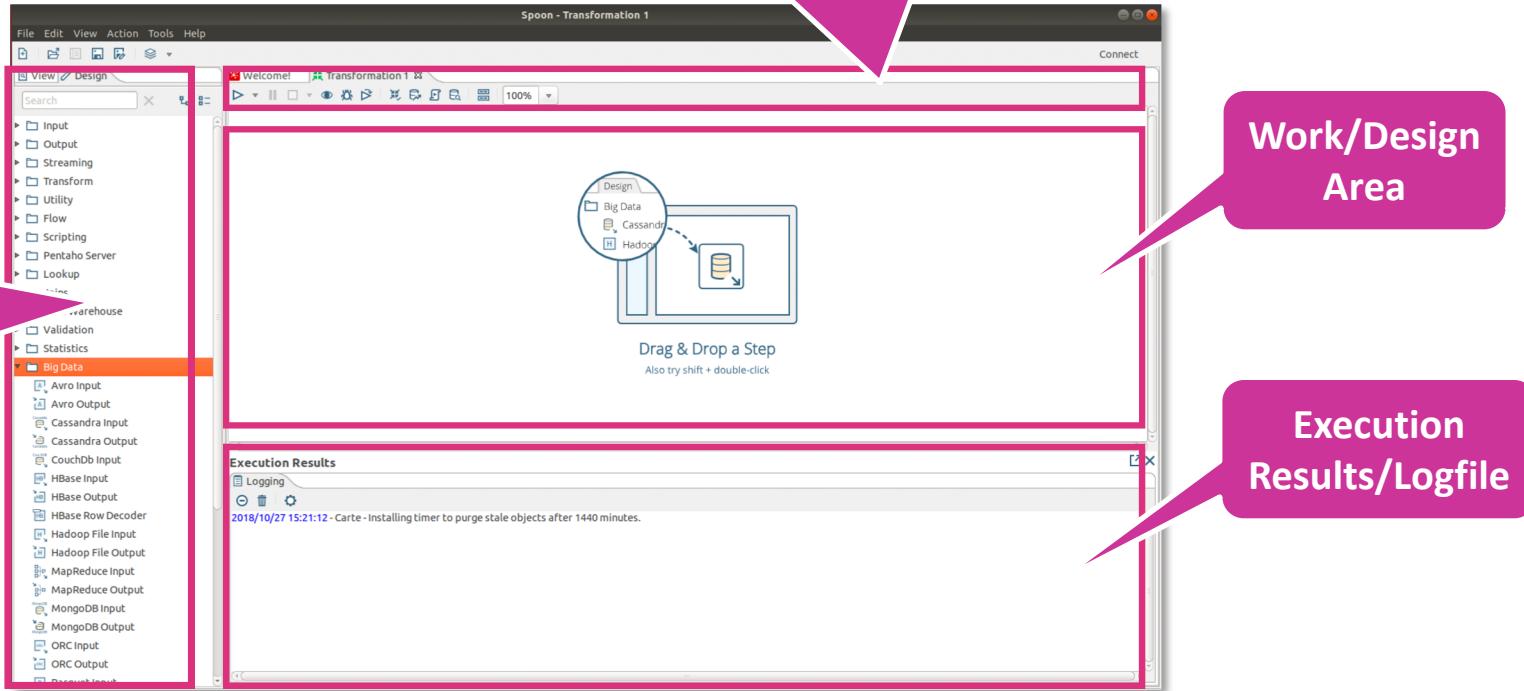
**UNIX:** /opt/pdi/spoon.sh

**Windows:** /some/directory/spoon.bat



# Spoon Interface

## 3. Start Pentaho Data Integration



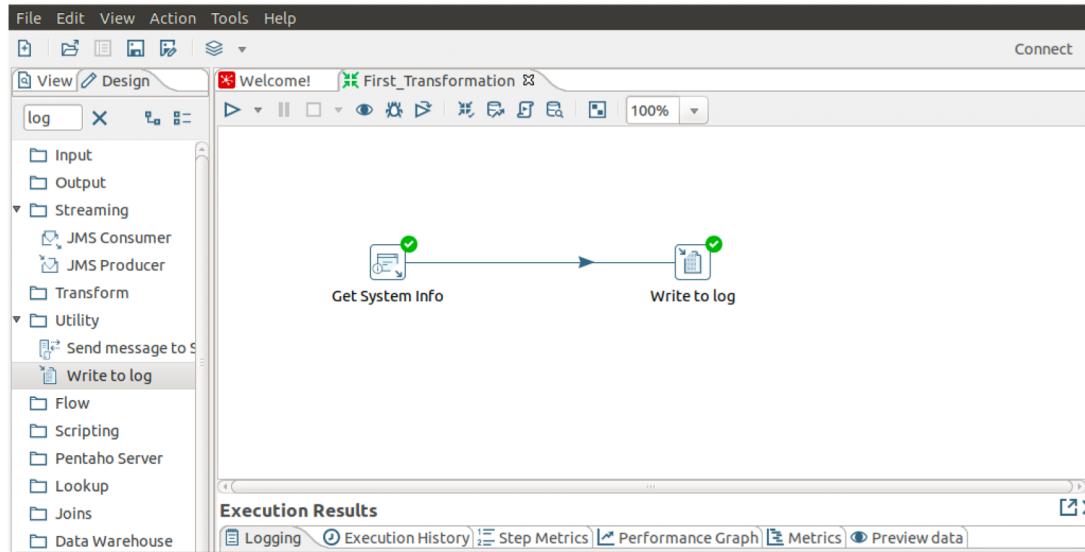
# Exercises Preparation II

Some naive ETL examples of using Pentaho Data Integration



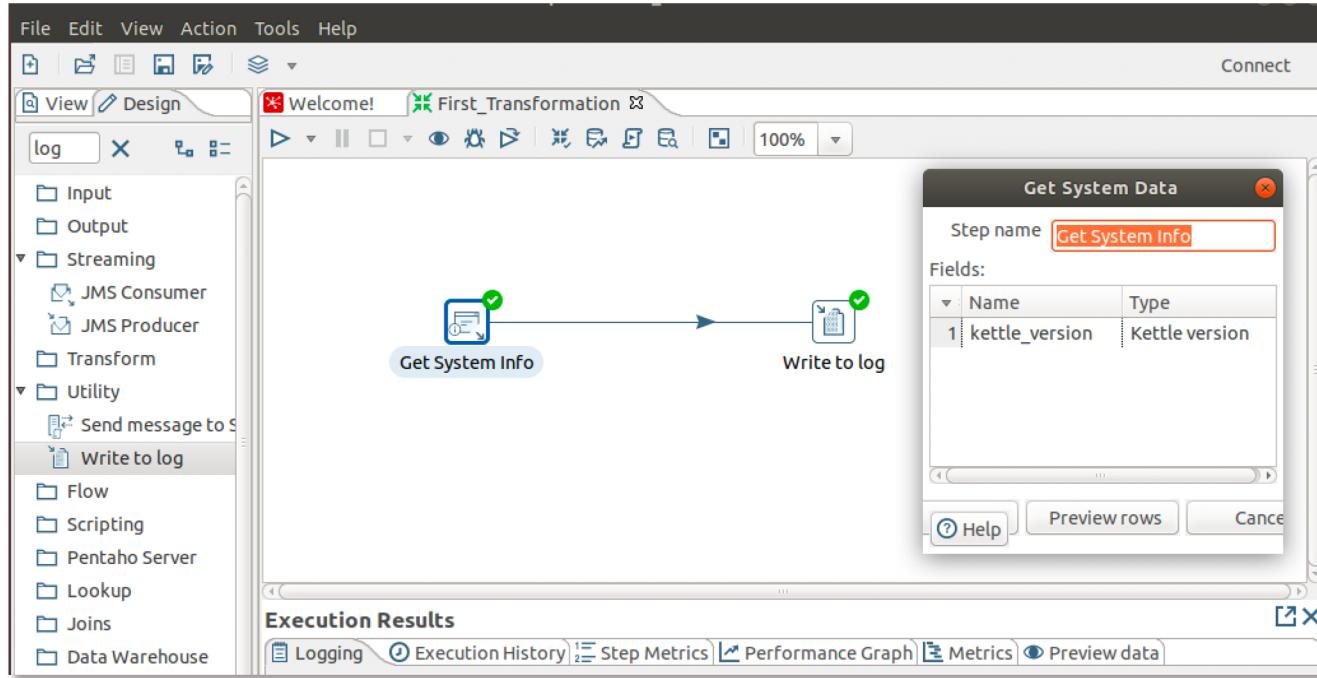
# PDI Basics/Examples – First Transformation

1. Create new **Transformation** (Click: *File* → *New* → *Transformation*)
2. Drag&Drop step „Get System Info“ and „Write To Log“ into **Work/Design Area** and connect both steps:



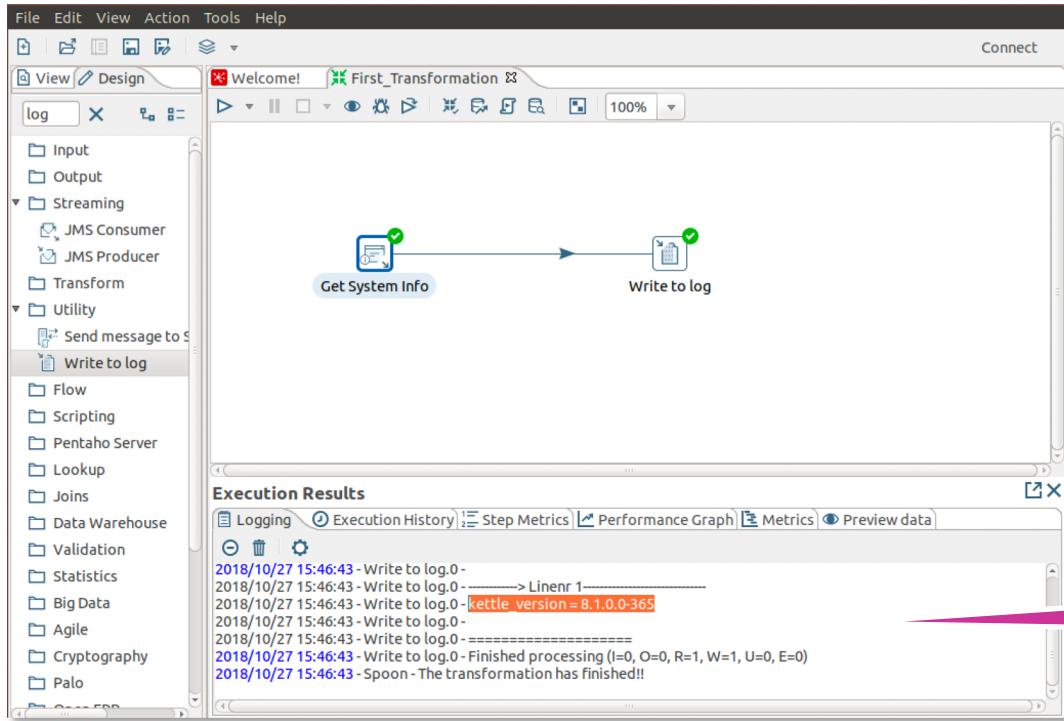
# PDI Basics/Examples – First Transformation

3. Double Click „Get System Info“ step to configure step accordingly to kettle (PDI) version info:



# PDI Basics/Examples – First Transformation

4. Save Transformation (*First\_Transformation*). Press Run Button. See Results:



- results of „Get System Info“ step are redirected to „Write To Log“ step
- „Write to Log“ step will write results to execution log
- **transformations** are all about actual **data transformation** and data flow

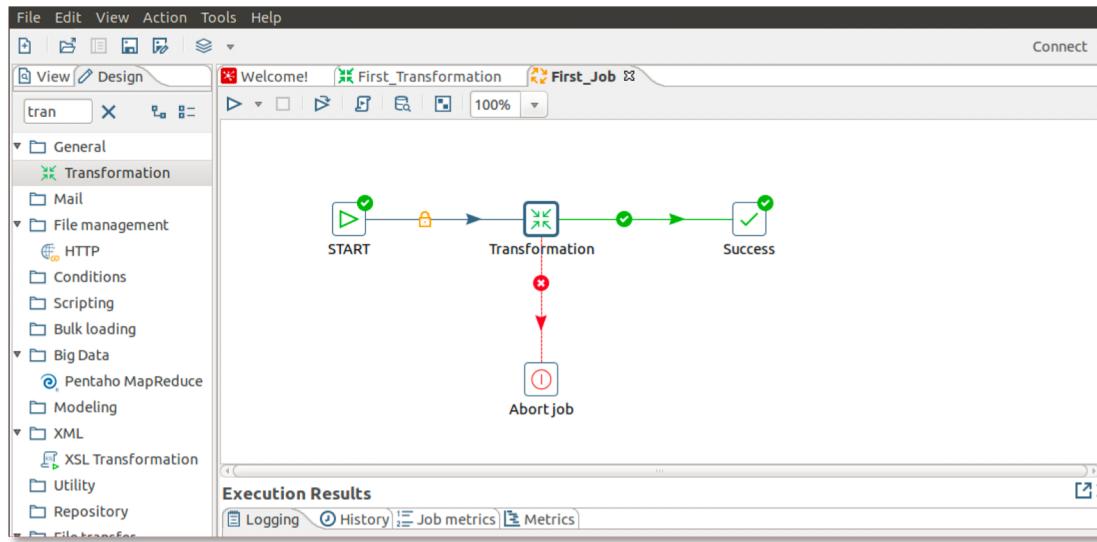
Results



# PDI Basics/Examples – First Job

1. Create New **Job** (Click: *File* → *New* → *Job*)

2. Drag&Drop step „**START**“, „**Transformation**“, „**Abort Job**“ and „**Success**“ steps into **Work/Design Area** and connect all steps accordingly:



- Start of each Job



- End of a Job (if not successful)



- End of a Job (if successful)

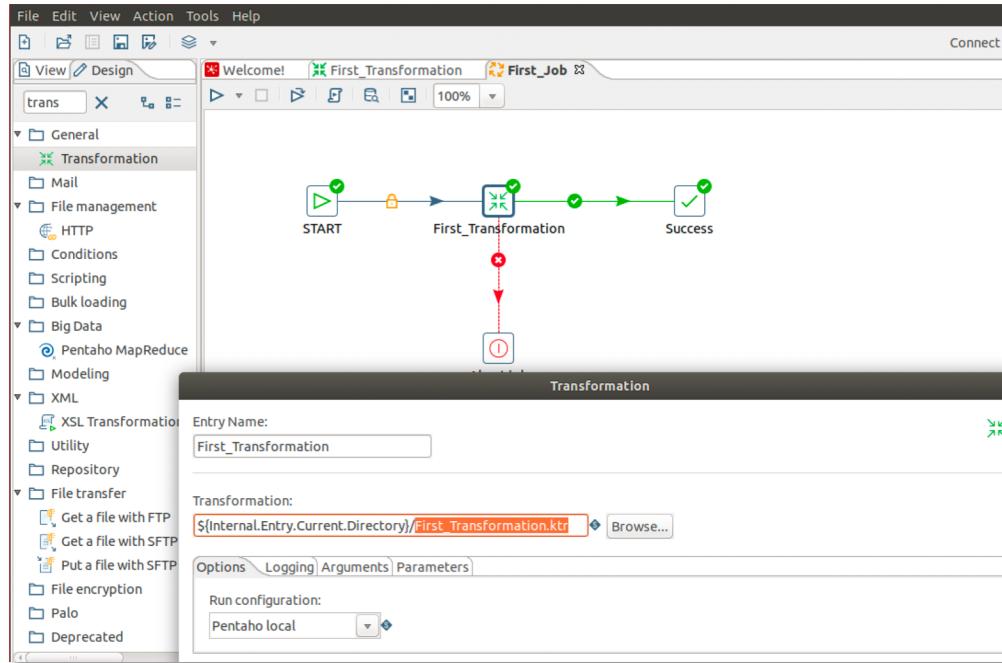


- Includes a Transformation



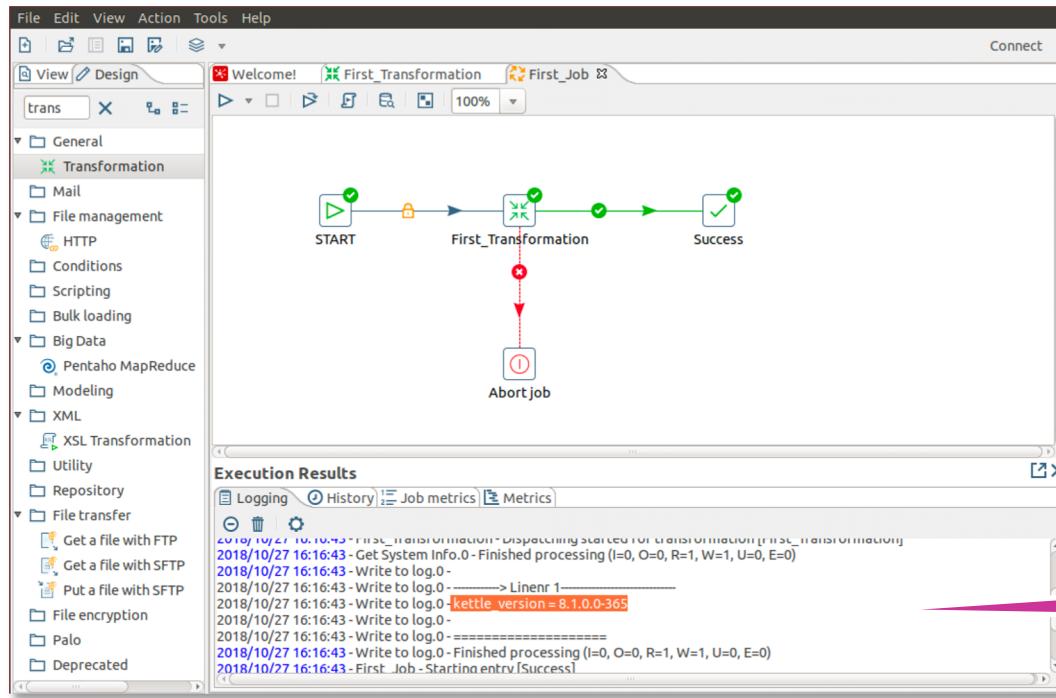
# PDI Basics/Examples – First Job

3. Include previously created Transformation (First\_Transformation) by double click on step „Transformation“ and including „First\_Transformation.ktr“:



# PDI Basics/Examples – First Job

4. Save Job (*First\_Job*). Press Run Button. See Results:



- Job will execute previously created transformation (*First\_Transformation.ktr*)
- Output of transformation will be appended to execution log of Job
- **jobs** are all about **workflows** of multiple **transformations** and **jobs**

Results

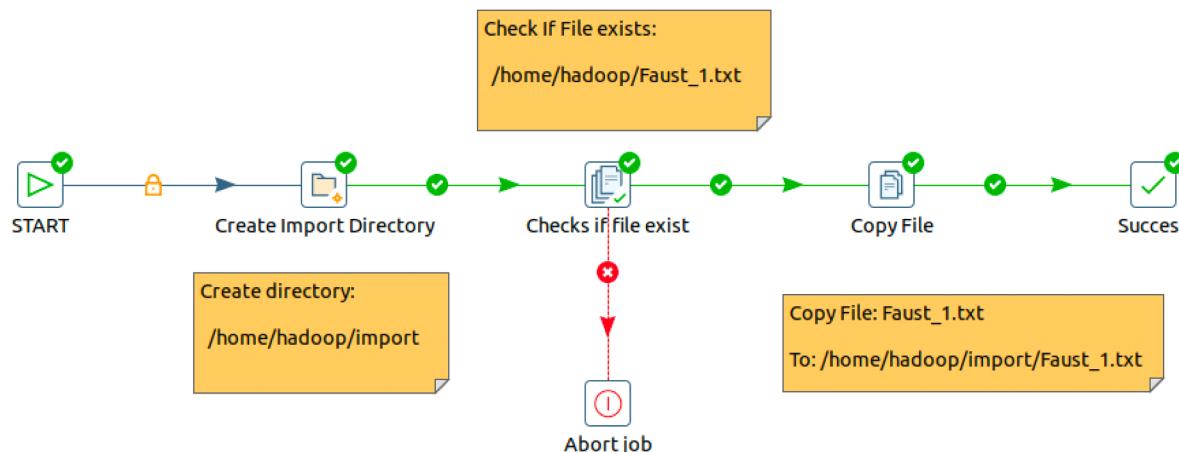


# PDI Basics/Examples – Local Filesystem

1. Create New **Job** *Local\_FileSystem.kjb* using steps:

- „Create a Folder“,
- „Check if files exist“ and „Copy File“

...to check whether `/home/hadoop/Faust_1.txt` exists and either abort (if not) or copy it to import directory `/home/hadoop/import/Faust_1.txt`



# PDI Basics/Examples – HDFS

## 1. Start HDFS and Yarn:

```
start-dfs.sh  
start-yarn.sh
```

## 2. Create New **Job** *HDFS\_Filesystem.kjb* using steps:

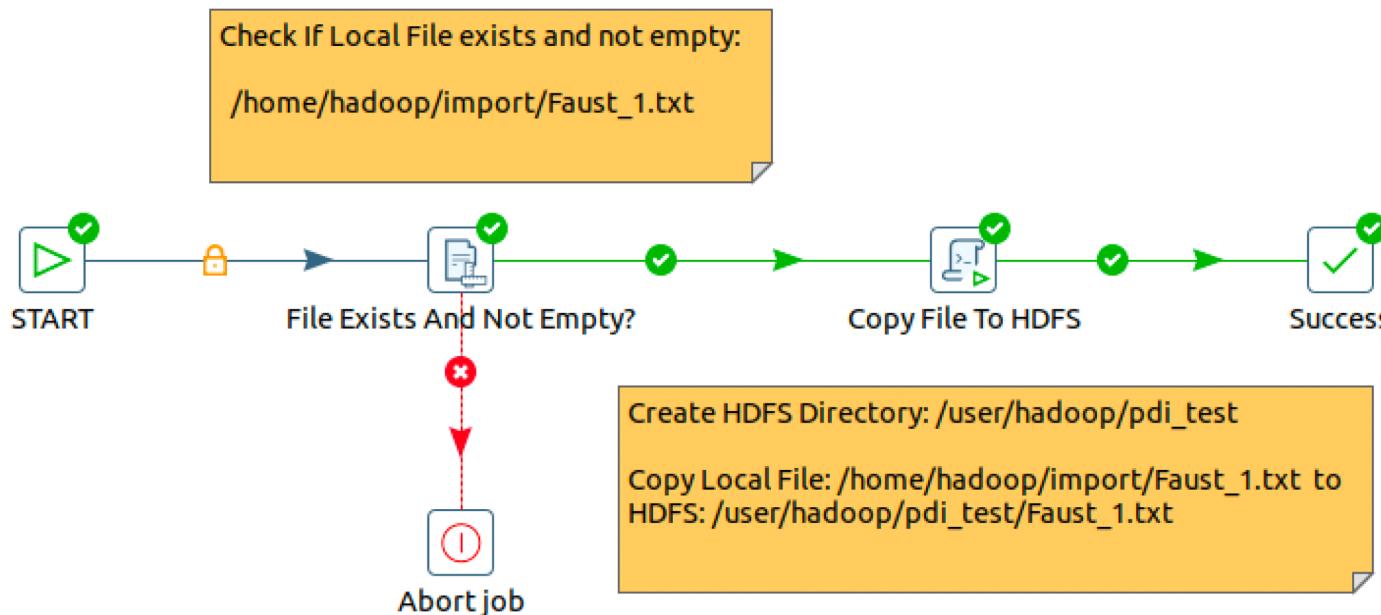
„Evaluate file metrics“ – to check whether file exists and is not empty

„Shell“ – to execute bash commands (e.g. `hadoop fs -put ...` in our case)



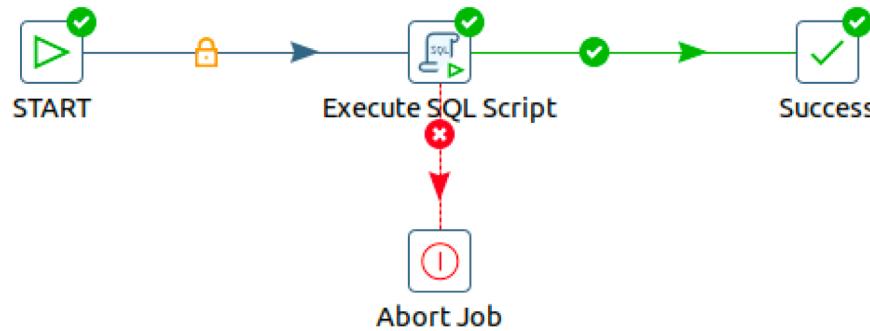
# PDI Basics/Examples – HDFS

3. Save and run job:



# PDI Basics/Examples – Hive

1. Create New **Job** *Hive\_SQL.kjb* using step „*Execute SQL Script*“:

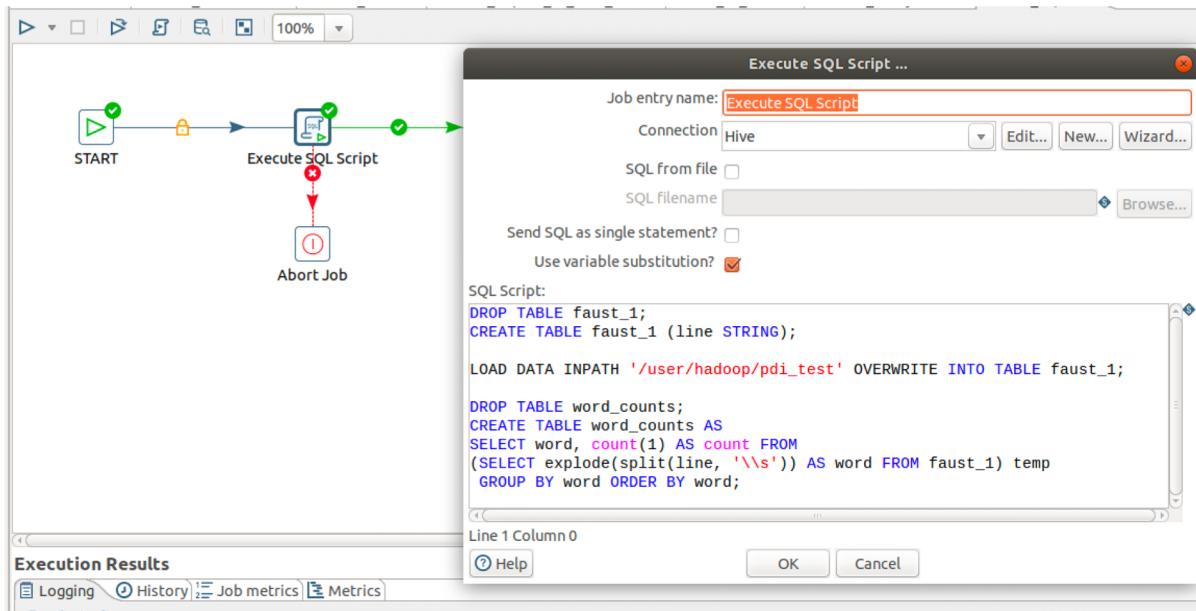


2. Start HiveServer2 (to enable job to connect to Hive):

```
/home/hadoop/hive/bin/hiveserver2
```

# PDI Basics/Examples – Hive

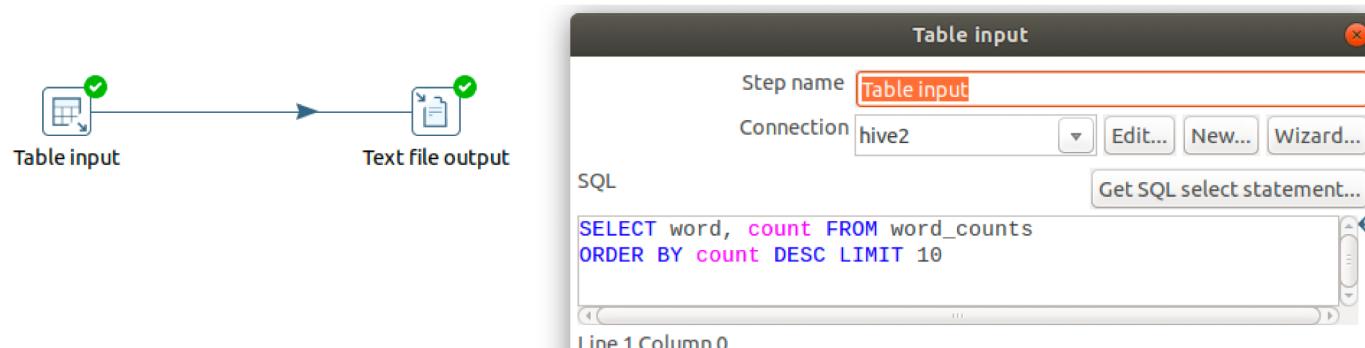
3. Add HiveQL statement to step „Execute SQL Script“ (word count example of previous lecture) and run job:



- Job will create table `faust_1`
- Load file `Faust_1.txt` from previous step into table `faust_1`
- Calculate word counts and save results to table `word_counts`

# PDI Basics/Examples – Parameters

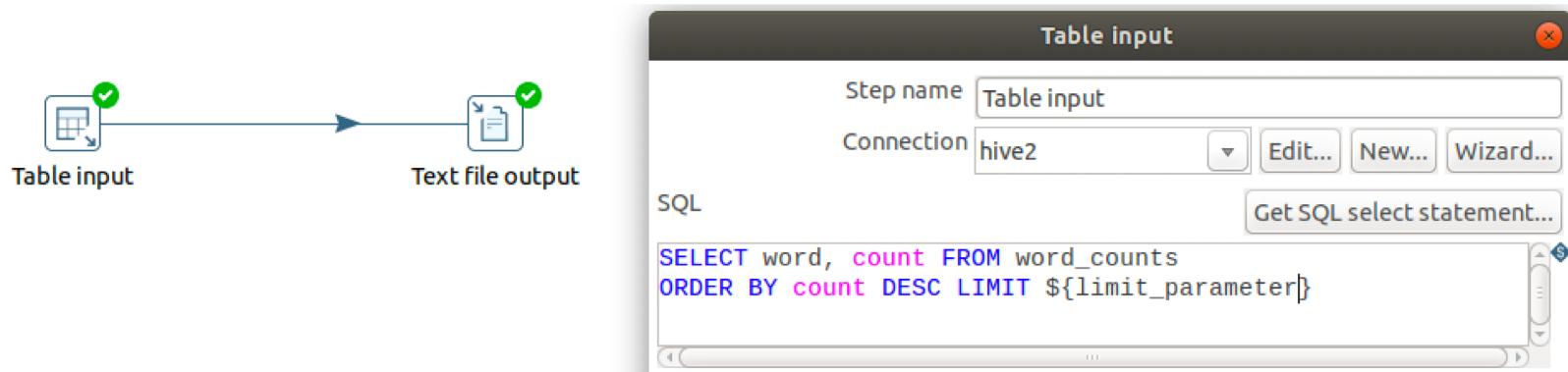
1. Create New **Transformation** `WordCount_Parameter.ktr` using steps „Table Input“ and „Text File output“:



- „Table Input“ reads word count results from table `word_counts` from previous step
- „Text file output“ saves result (Top 10 word counts) of „Table input“ query to csv file on local filesystem

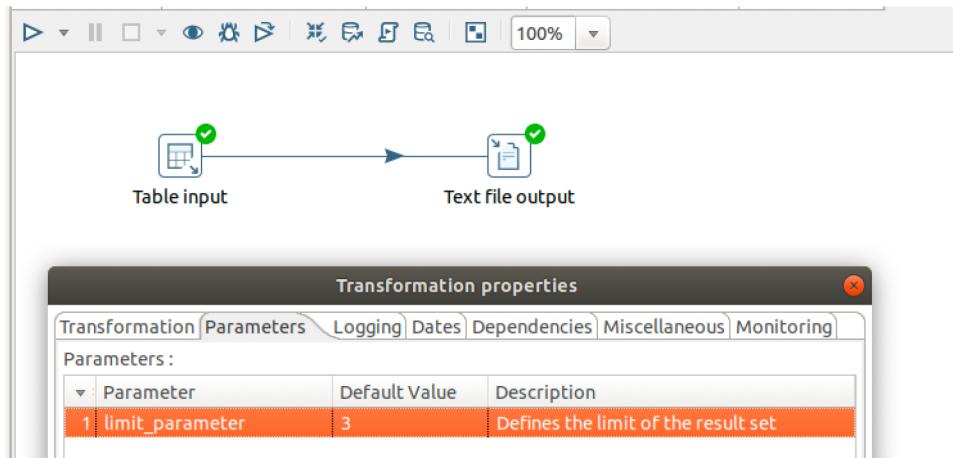
# PDI Basics/Examples – Parameters

2. Stop! **Parameters?** We will make use of a parameter to decide how many results we want to have inside the result set (local CSV file):



# PDI Basics/Examples – Parameters

3. Parameters can be set using Edit → Settings → Parameters



4. Run job. Output CSV will only have 3 entries now.

```
hadoop@marcel-VirtualBox:~$ cat export/word_count.csv
word,count
        ,1603
und    ,509
die   ,463
```

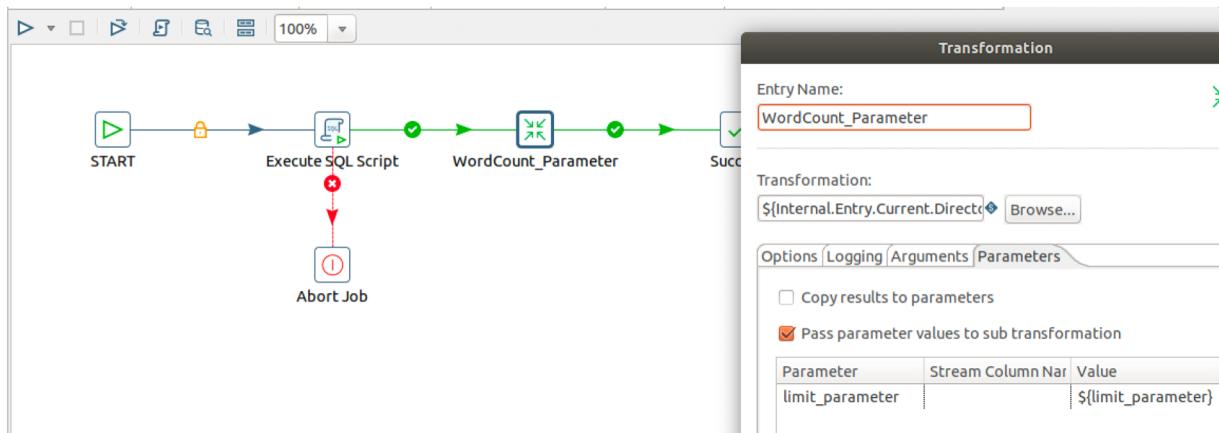


# PDI Basics/Examples – Parameters

5. Parameters can also be set within `kettle.properties` file:

```
hadoop@marcel-VirtualBox:~$ vi .kettle/kettle.properties
limit_parameter=7
```

6. Parameters can also be passed from Job To Job/Transformation, e.g. by extending previous example job „`Hive_SQL.kjb`“ by transformation „`WordCount Parameter.ktr`“ and passing `kettle.properties` parameter `limit_parameter` to it:



# PDI Basics/Examples – Execute Jobs Using Kitchen

1. It's nice to run jobs locally during development, but how to run them on a remote server (e.g. in a productive manner)? This is easily achieved using kitchen.sh:

```
/opt/pdi/kitchen.sh -file=/home/marcel/HDFS_Filesystem.kjb

[...]
2018/10/27 20:41:05 - HDFS_Filesystem - Start of job execution
log4j:ERROR No output stream or file set for the appender named [pdi-execution-appender].
2018/10/27 20:41:05 - HDFS_Filesystem - Starting entry [File Exists And Not Empty?]
2018/10/27 20:41:05 - HDFS_Filesystem - Starting entry [Copy File To HDFS]
2018/10/27 20:41:05 - Copy File To HDFS - Running on platform : Linux
2018/10/27 20:41:09 - HDFS_Filesystem - Starting entry [Success]
2018/10/27 20:41:09 - Carte - Installing timer to purge stale objects after 1440 minutes.
2018/10/27 20:41:09 - HDFS_Filesystem - Finished job entry [Success] (result=[true])
2018/10/27 20:41:09 - HDFS_Filesystem - Finished job entry [Copy File To HDFS] (result=[true])
2018/10/27 20:41:09 - HDFS_Filesystem - Finished job entry [File Exists And Not Empty?] (result
=[true])
2018/10/27 20:41:09 - HDFS_Filesystem - Job execution finished
2018/10/27 20:41:09 - Kitchen - Finished!
2018/10/27 20:41:09 - Kitchen - Start=2018/10/27 20:41:00.243, Stop=2018/10/27 20:41:09.265
2018/10/27 20:41:09 - Kitchen - Processing ended after 9 seconds.
```



# Exercises Preparation III

A simple ETL examples of how to use an ETL Workflow tool (PDI) to integrate IMDb data



# PDI IMDb Import

Previous examples have been nice to get a basic understanding of how PDI works, but not in a productive manner, as:

- There haven't been any complex workflows
- An appropriate usage of parameters
- Dependencies
- ...

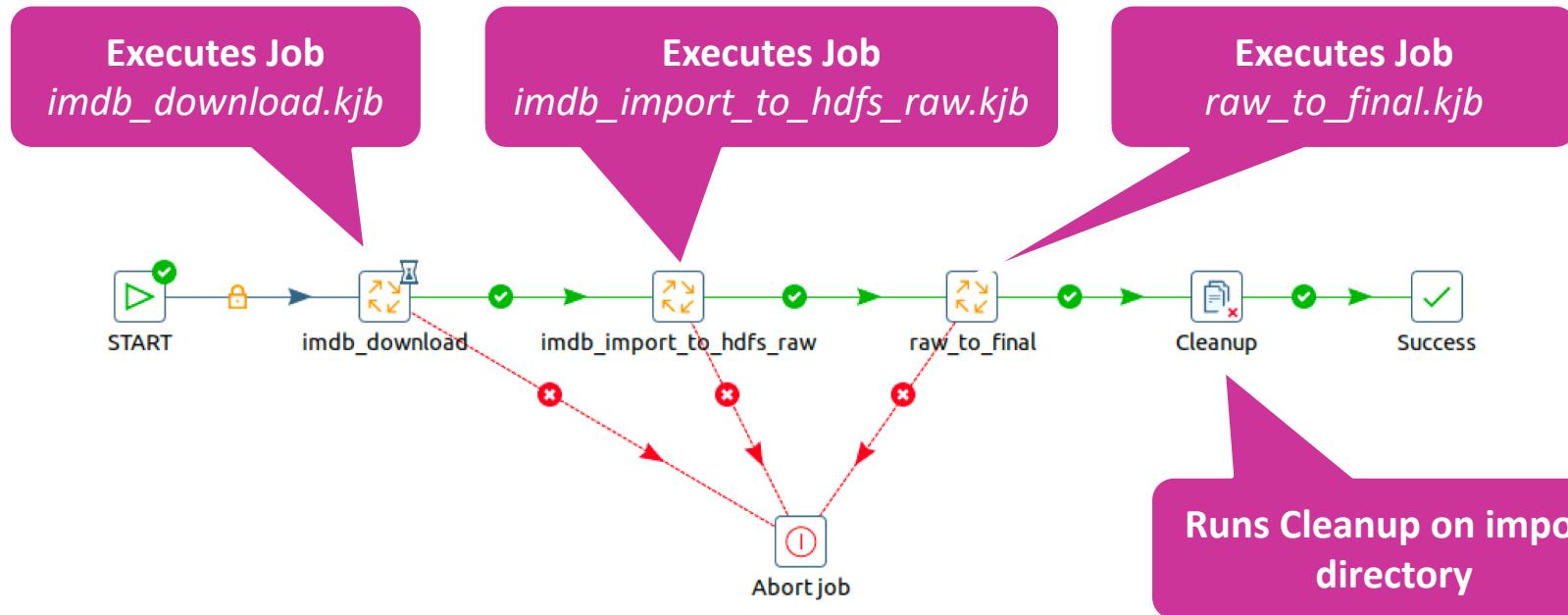
That's why we take a look at a real world example: **IMDB data**

And describe a complete workflow which is able to run each day:

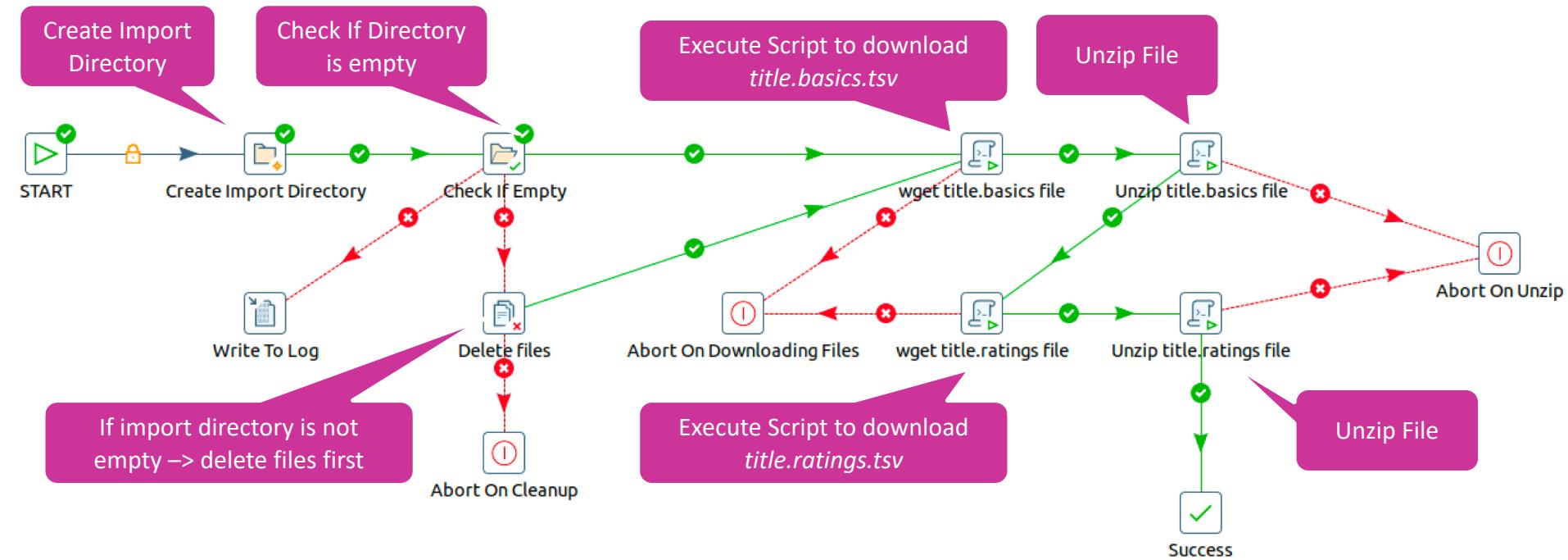
- **Download of IMDb data** to local filesystem
- Move Data to HDFS (raw directory/layer )
- Create Hive tables for **Raw Layer**
- Create and fill **Final Layer** (Hive tables) by raw layer tables, applying business rules and using dynamic partitioning



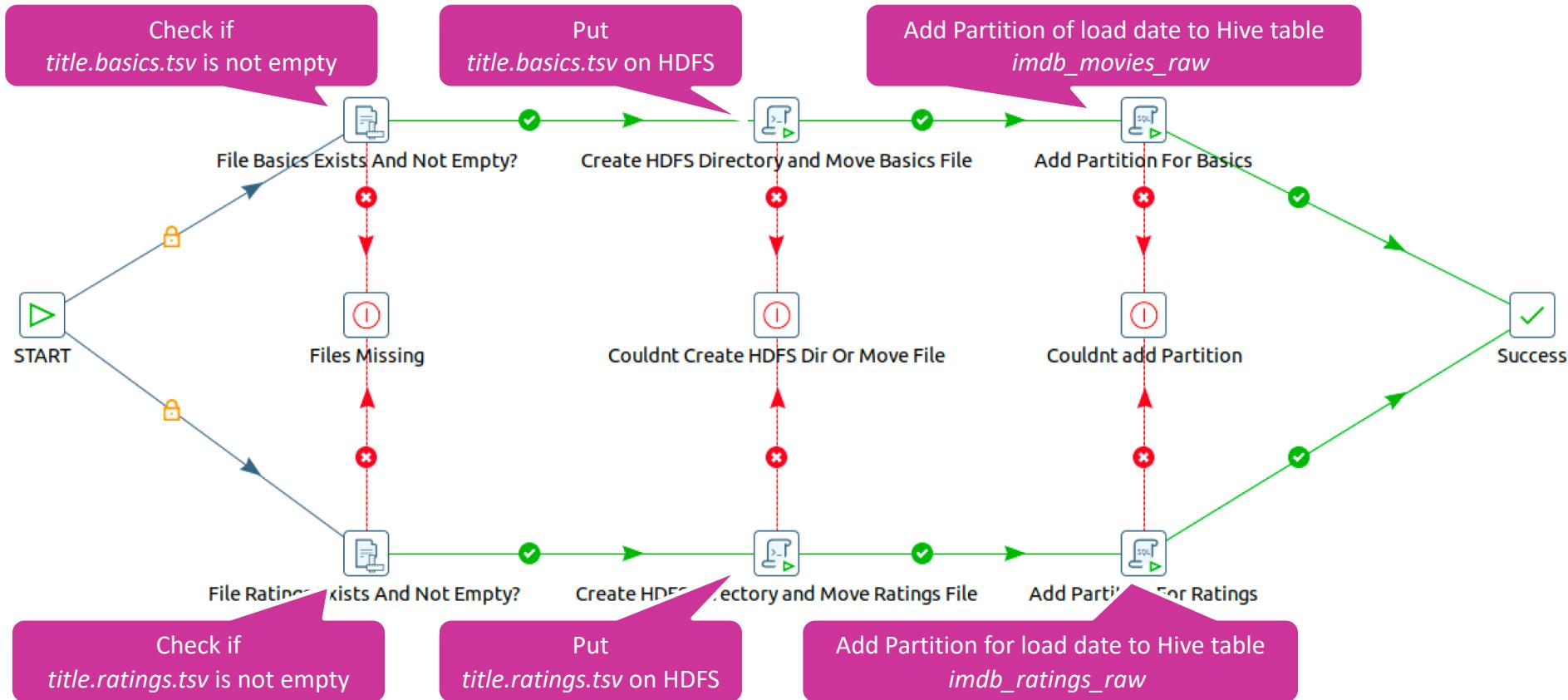
# PDI IMDb Import – Main Job



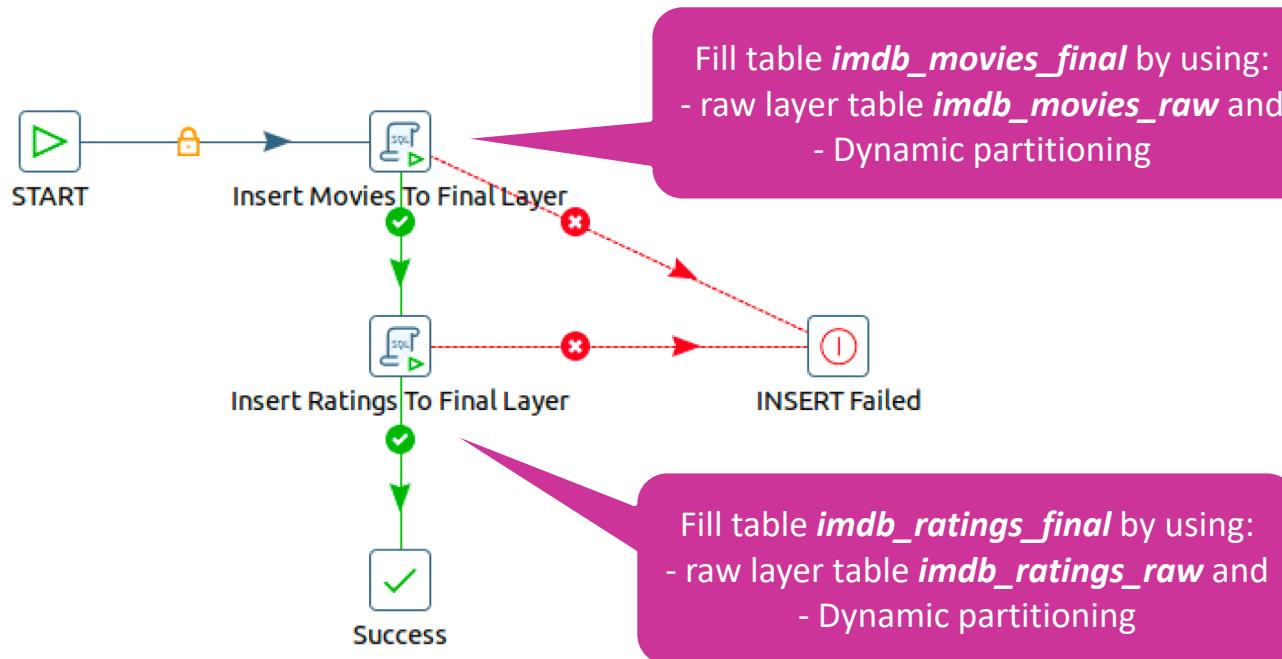
# PDI IMDb Import – *imdb\_download.kjb*



# PDI IMDb Import – *imdb\_import\_to\_hdfs\_raw.kjb*



# PDI IMDb Import – *raw\_to\_final.kjb*



# PDI IMDb Import – Execution

1. Execute using Spoon or kitchen.sh:

```
/opt/pdi/kitchen.sh -file=/home/hadoop/imdb_main.kjb -param:load_day=29  
-param:load_month=10 -param:load_year=2018
```

A Job like this could be scheduled by cron to run on a daily basis, receiving parameters:

- load\_day
- load\_month
- load\_year

... from cron job





# Exercises

Use Pentaho Data Integration to solve exercises  
based on IMDb data



# Pentaho Data Integration Exercises – IMDB

1. Execute Tasks of previous HandsOn Slides
2. Use PDI and previous **Jobs&Transformations** to do following changes:
  - a) **Extend job** *imdb\_download.kjb* to also download ***name.basics.tsv.gz***
  - b) **Extend job** *imdb\_import\_to\_hdfs\_raw.kjb* to also import ***name.basics.tsv*** to HDFS raw layer.
  - c) **Create Hive table** *imdb\_actors\_raw* for ***name.basics.tsv*** in raw layer.  
Table should be partitioned by year, month and day of load date.
  - d) **Create table** *imdb\_actors\_final* and **extend job** *raw\_to\_final.kjb* to also fill table *imdb\_actors\_final* using:
    - data of table *imdb\_actors\_raw* and
    - **partition table** by column *partition\_is\_alive* containing „alive“ or „dead“, wether the actor is alive or dead.



# Pentaho Data Integration Exercises – IMDB

3. Run main workflow job **imdb\_main.kjb** using:

```
/opt/pdi/kitchen.sh -file=/home/hadoop/imdb_main.kjb -param:load_day=29  
-param:load_month=10 -param:load_year=2018
```

