Open for Innovation

# KNIME

# Re-engineering IoT Legacy Analytics Solutions with Big Data

Rosaria Silipo & Bernd Wiswedel
KNIME.com

Rosaria.Silipo@knime.com
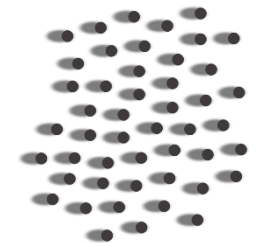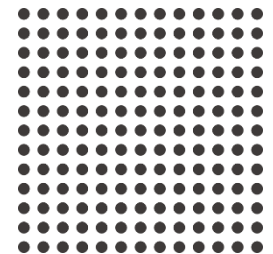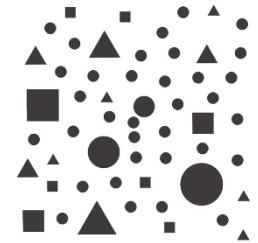
# Variety, Volume, Velocity

Variety:

- integrating heterogeneous data (and tools)

Volume:

- from small files…
- …to distributed data repositories (Hadoop)
- bring the tools to the data

Velocity:

- from distributing computationally heavy computations…
- …to real time scoring of millions of records/sec.

# Every Minute...

# IoT



THE INTERNET OF THINGS

AN EXPLOSION OF CONNECTED POSSIBILITY

# The IoT Legacy

KNIME
Open for Innovation

# Energy Usage Prediction from Smart Meters Data

- Read Smart Meter Energy Data (**176 millions rows**)
- Clean Up and Aggregate total Energy Usage by hour, week, day, month, year

  Workflow 1
- Calculate Behavioral Measures for each Smart Meter

- Cluster Smart Meters with Similar Behavior (k-Means)

  Workflow 2

- Predict Energy Usage in Clustered Smart Meters (Auto-Regressive Time Series Prediction)

  Workflow 3

Open for Innovation
KNIME

# Workflow 1: PrepareData



This workflow reads Ireland's electricity data, converts the dates from the proprietary format into datetime values, and groups kW values by:
- day
- hour
- intra-day times
- month
- year
- week

It also aggregates average and % values for the k-Means procedure

**Read all Data**

Read 6 files for a total of 176 Mio Rows

**String to datetime**

convert proprietary date format into datetime values

**Daily, Monthly, Yearly, Weekly**

kW usage by meter ID by day, month, week, year
The top port also offers:
average kW usage daily, monthly, weekly, yearly by meter ID

**Hourly, Intra-day**

kW usage by meter ID by hour and intra-day times
The top port also offers:
average kW usage hourly, for each intra-day time by meter ID

**Joiner**

Node 50

**% values**

intra-day and intra-week kW % usage by meter ID

**Write to server**

Node 99

**Write to CSV**

Node 100

# ~ 2 days

# Big Data Options

8

# Big Data Support

- KNIME Big Data Access Nodes
  - preconfigured connectors
  - in database processing
- Big Data Platforms
  - HDFS, Hive, Impala, HP Vertica, Hortonworks, ParStream, Actian, MapR, any big data platform really!
- Spark MLlib integration (coming soon)
- Streaming Executor (coming soon)

# Hadoop Sandboxes

- Hortonworks:

  http://hortonworks.com/products/hortonworks-sandbox/

- Cloudera:

  http://www.cloudera.com/content/cloudera/en/downloads/quickstart_vms.html

- Virtual Box

  https://www.virtualbox.org/

- VMWare Player

  http://www.vmware.com/

# … as easy as 1,2,3,… 4

**1**    **2**    **3**    **4**

| Access Big Data | Select Table | In-DB Processing | Into KNIME |
|---|---|---|---|

# 1. Database Connector

Generic Database Connector
- Can connect to any JDBC source
- Register new JDBC driver via preferences page

Open for Innovation
KNIME

# 1. Register JDBC Driver

Increase connection timeout for long running retrieval operations

Open KNIME and go to File -> Preferences

# 1. Dedicated Connectors

Dedicated pre-configured connectors

- – Bundling necessary JDBC drivers
- – Easy to use
- – DB specific behavior/capability

**BigData Connectors**

**Impala Connector**

works for most Hadoop HIVE installations, including **Hortonworks**

**Hive Connector**

Some dedicated connectors are part of the open source KNIME Analytics Platform, some belong to the commercial KNIME Big Data Extension

free

**Vertica Connector**

Open for Innovation

KNIME

# 2. Data Table Selection

**Connect to a big data platform:**
- Impala
- Hive
- parStream
to read the energy data

**Database Connector**

parStream platform

**Impala Connector**

Cloudera Amazon cluster
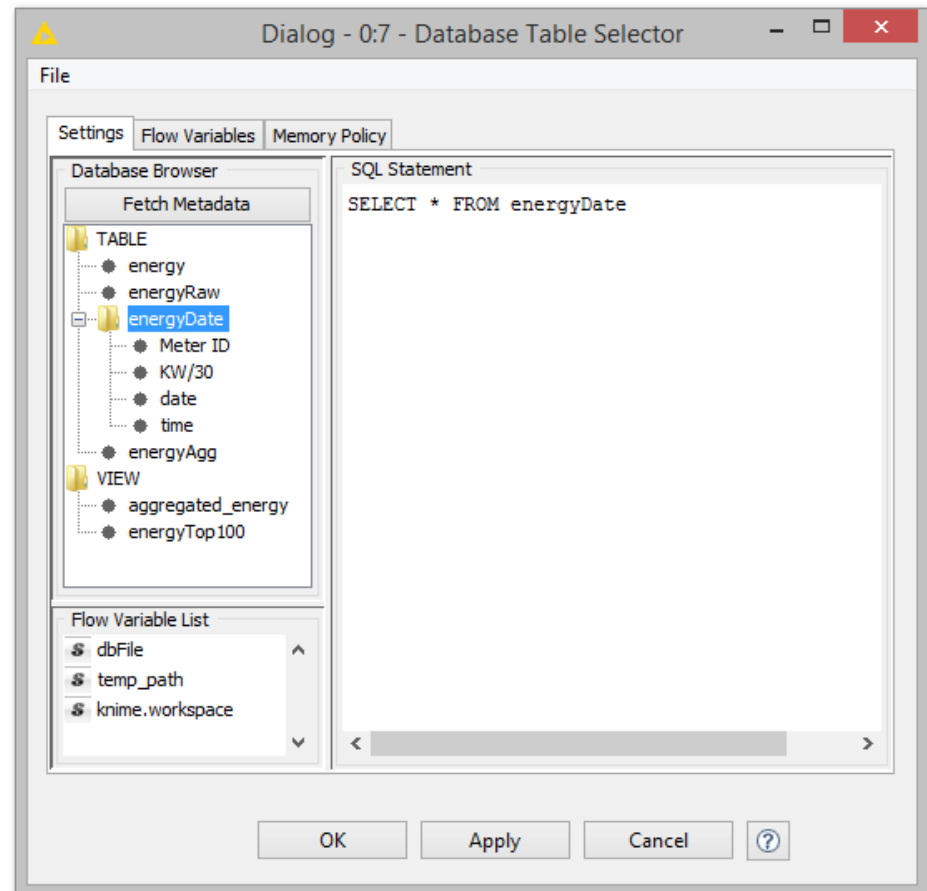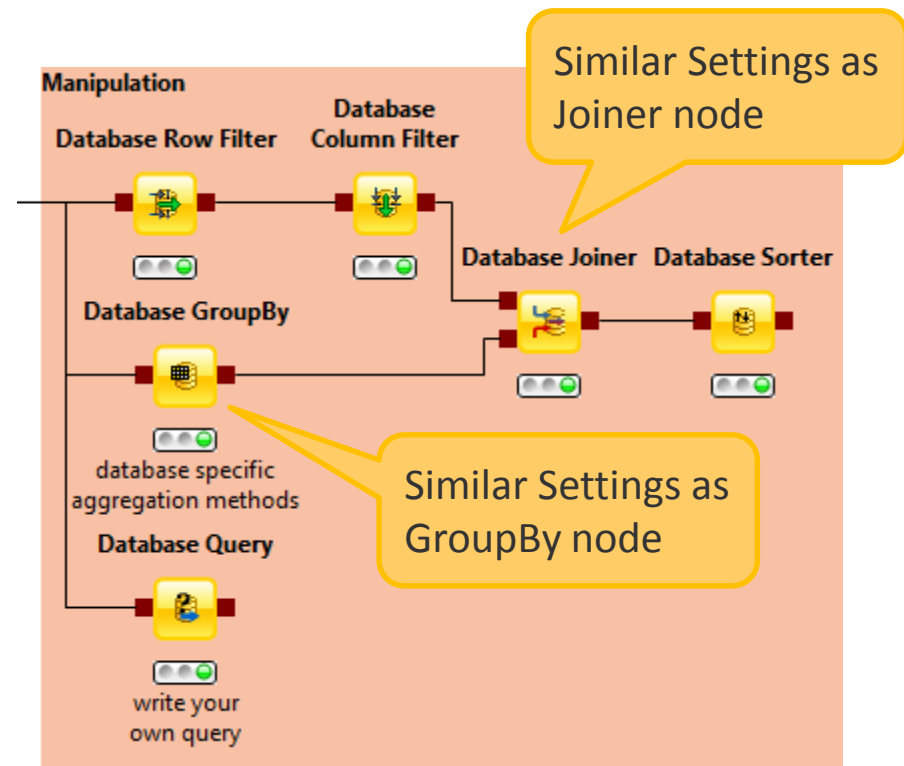
**Database Table Selector**

table "energy" with all energy measurements sampled every half an hour one year long almost 6000 meter IDs
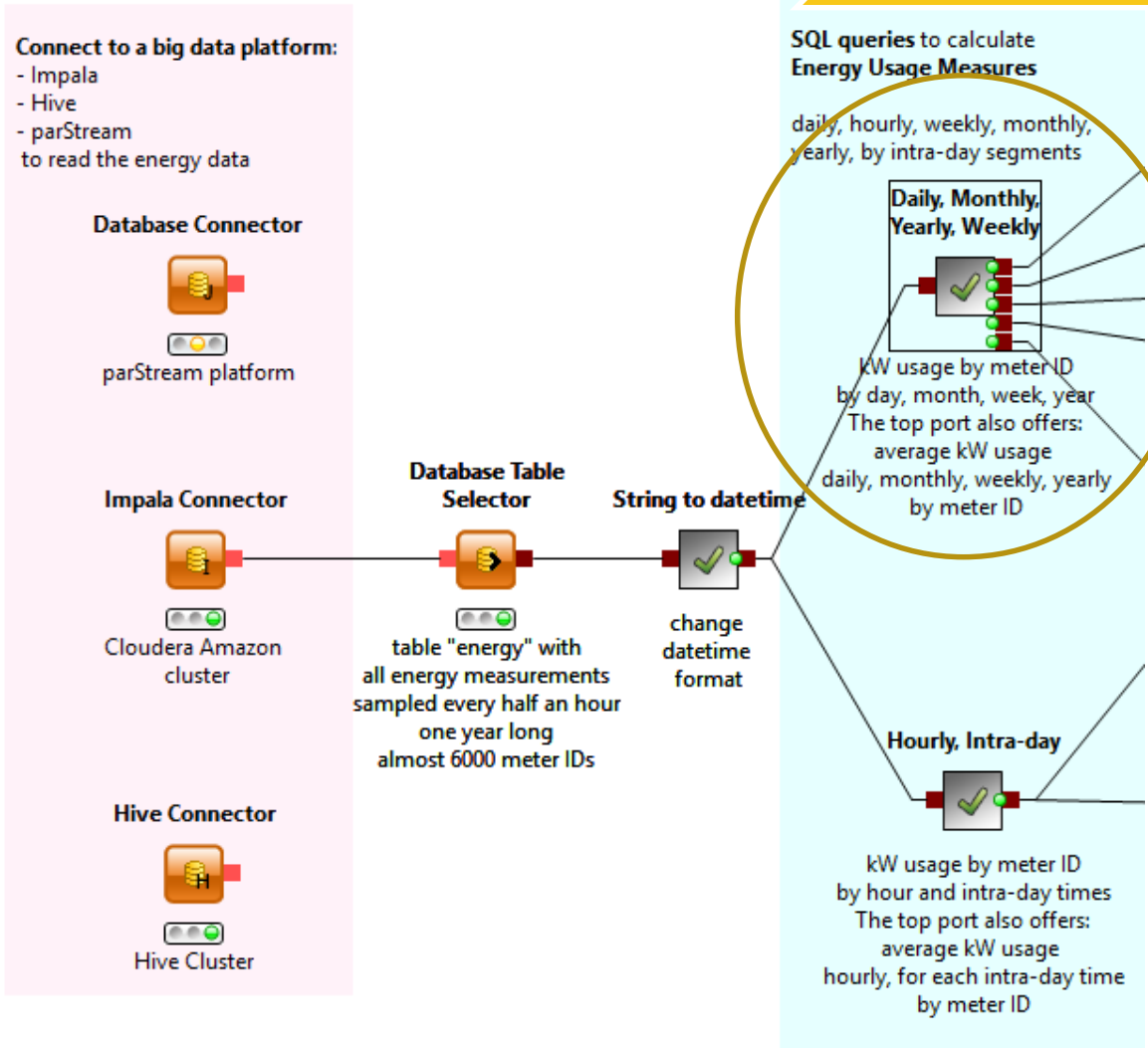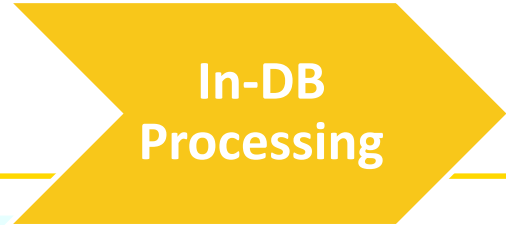
**Hive Connector**

Hive Cluster

Dialog - 0:7 - Database Table Selector

File

Settings | Flow Variables | Memory Policy

Database Browser

Fetch Metadata

TABLE
- energy
- energyRaw
- energyDate
  - Meter ID
  - KW/30
  - date
  - time
- energyAgg
VIEW
- aggregated_energy
- energyTop100

SQL Statement

SELECT * FROM energyDate

Flow Variable List
- *s* dbFile
- *s* temp_path
- *s* knime.workspace

OK    Apply    Cancel

Open for Innovation
KNIME

# 3. In-Database Processing

- Filter rows and columns

- Join tables/queries

- Sort your data

- Write your own query

- Aggregate* your data



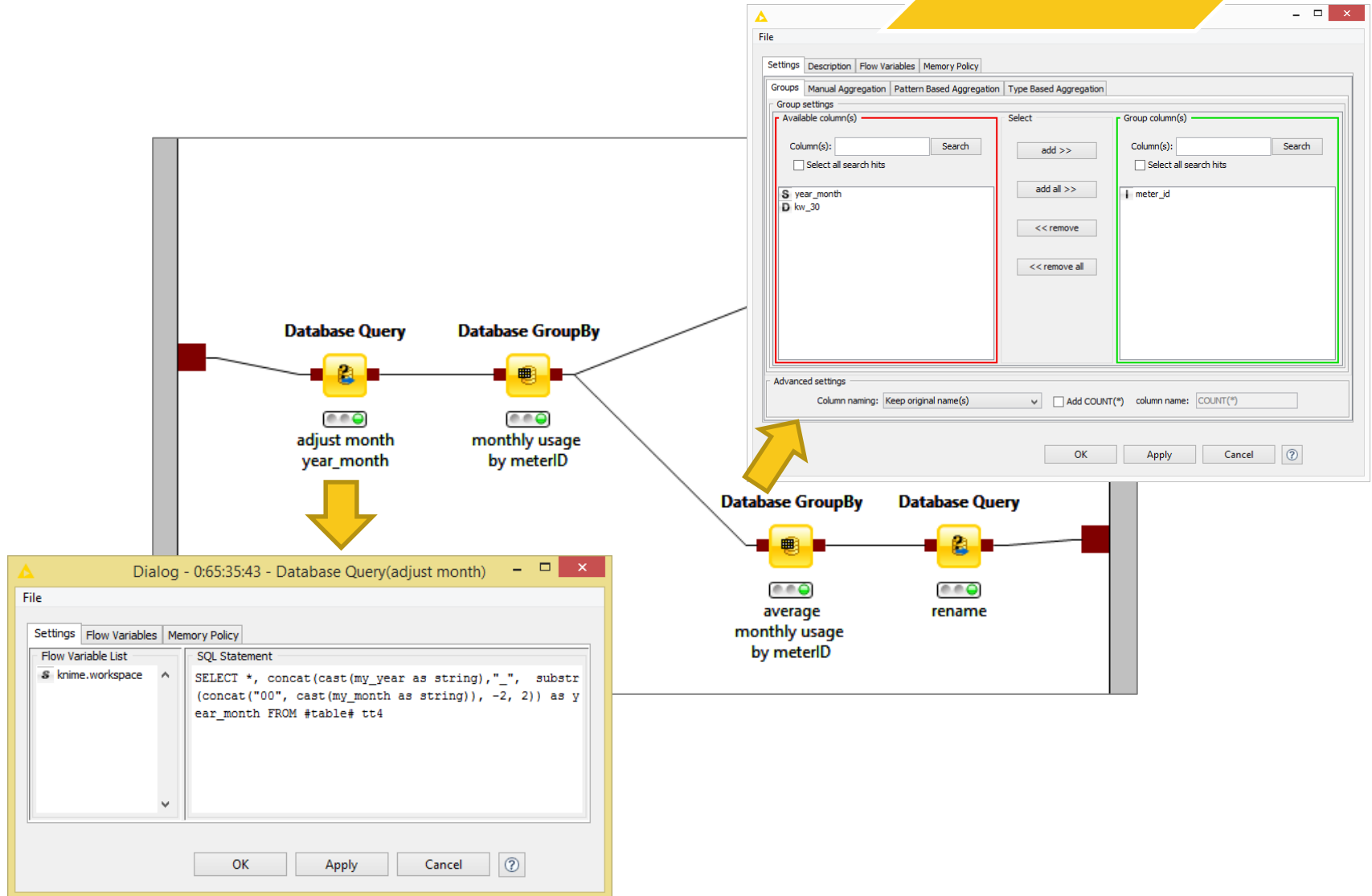Similar Settings as Joiner node

Similar Settings as GroupBy node

* Database GroupBy node exposes DB specific aggregation methods

# 3. Queries for average Measures

**Connect to a big data platform:**
- Impala
- Hive
- parStream
to read the energy data

**Database Connector**

parStream platform

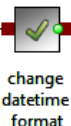**Impala Connector**

Cloudera Amazon cluster

**Database Table Selector**

table "energy" with all energy measurements sampled every half an hour one year long almost 6000 meter IDs

**String to datetime**

change datetime format

**Hive Connector**

Hive Cluster

**SQL queries** to calculate **Energy Usage Measures**

daily, hourly, weekly, monthly, yearly, by intra-day segments

**Daily, Monthly, Yearly, Weekly**

kW usage by meter ID by day, month, week, year The top port also offers: average kW usage daily, monthly, weekly, yearly by meter ID
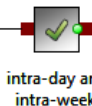
**Hourly, Intra-day**

kW usage by meter ID by hour and intra-day times The top port also offers: average kW usage hourly, for each intra-day time by meter ID
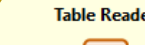
Open for Innovation
KNIME

# 3. Average Monthly Values

This workflow reads Ireland's electricity data,
converts the dates from the proprietary format into datetime values,
and groups kW values by, day, hour, intra-day times, month, year, week

It also aggregates average and % values for the k-Means procedure

**3**

**Connect to a big data platform:**
- Impala
- Hive
- parStream
to read the energy data

**Database Connector**

**1**

arStream platform

**2**

**SQL queries to calculate
Energy Usage Measures**

daily, hourly, weekly, monthly,
yearly, by intra-day segments

**Daily, Monthly,
Yearly, Weekly**

**Table Reader**

monthly
**Database Connection
Table Reader**

daily

**Database Connection
Table Reader**

**4**

yearly
**Database Connection
Table Reader**

**Impala Connector**

Cloudera Amazon
cluster

**Database Table
Selector**

table "energy" with
all energy measurements
sampled every half an hour
one year long
almost 6000 meter IDs

**String to datetime**

change
datetime
format

kW usage by meter ID
by day, month, week, year
The top port also offers:
average kW usage
daily, monthly, weekly, yearly
by meter ID

weekly
**Database Connection
Table Reader**

**Database Joiner**

join time series

**% values**

intra-day and
intra-week
kW % usage
by meter ID

result

**<30 min**

**Hive Connector**

Hive Cluster

**Hourly, Intra-day**

kW usage by meter ID
by hour and intra-day times
The top port also offers:
average kW usage
hourly, for each intra-day time
by meter ID

atabase Con..ction
able Read.

hourly

Run the SQL queries and
**retrieve final data**

Open for Innovation
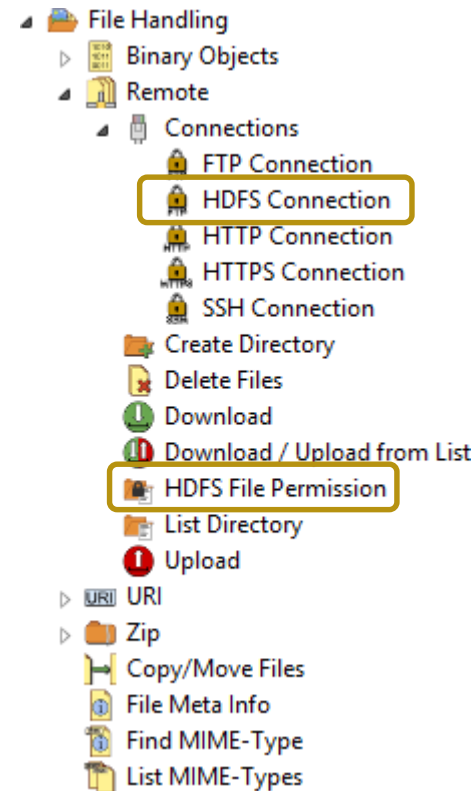KNIME

# New Big Data Platform?

# Other Useful Database Nodes

- Drop table
  - missing table handling
  - cascade option
- Execute any SQL statement
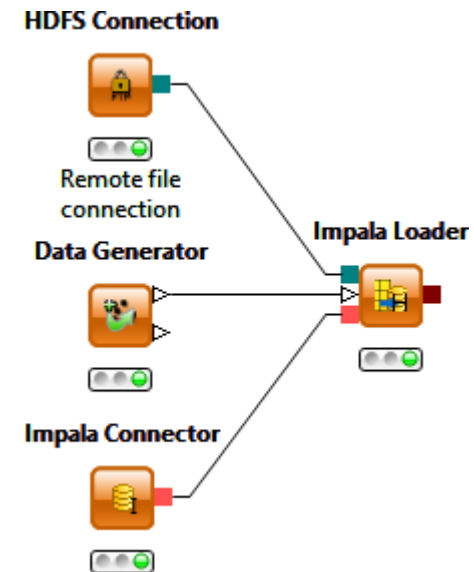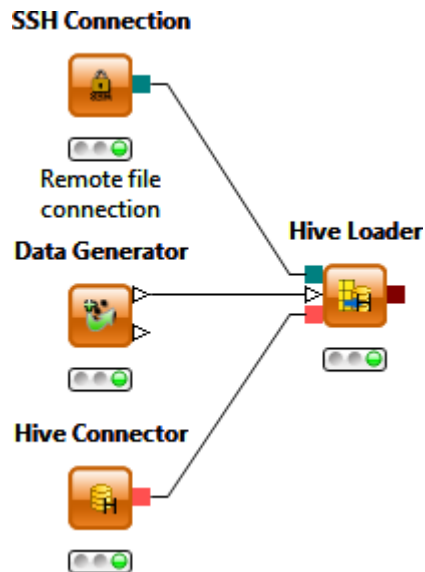- Manipulate existing queries

# HDFS File Handling

- KNIME & Extensions -> KNIME File Handling Nodes

- HDFS Connection and HDFS File Permission nodes

# Hive/Impala Loader
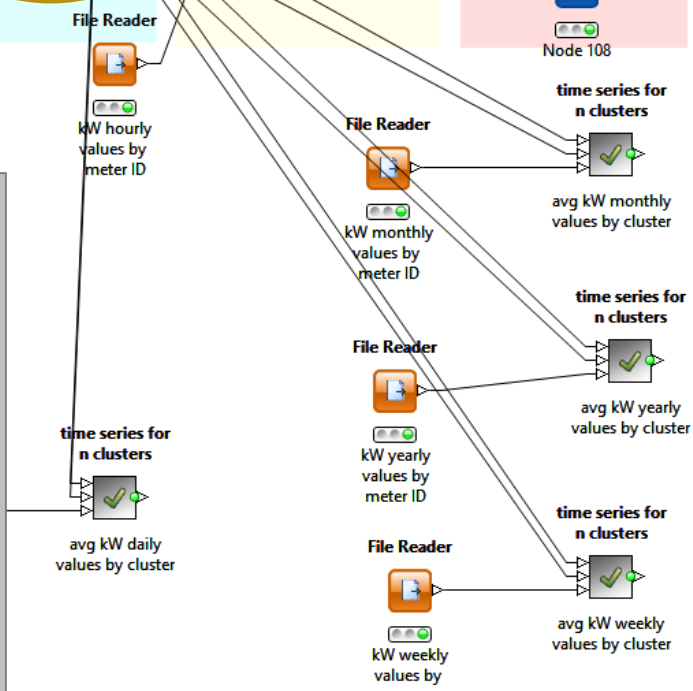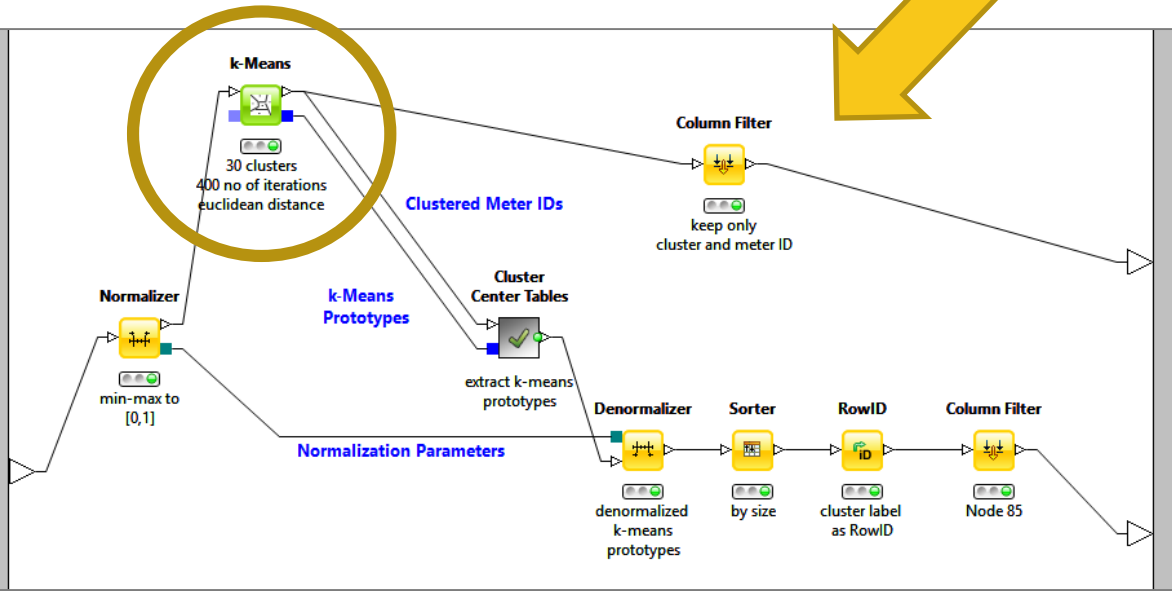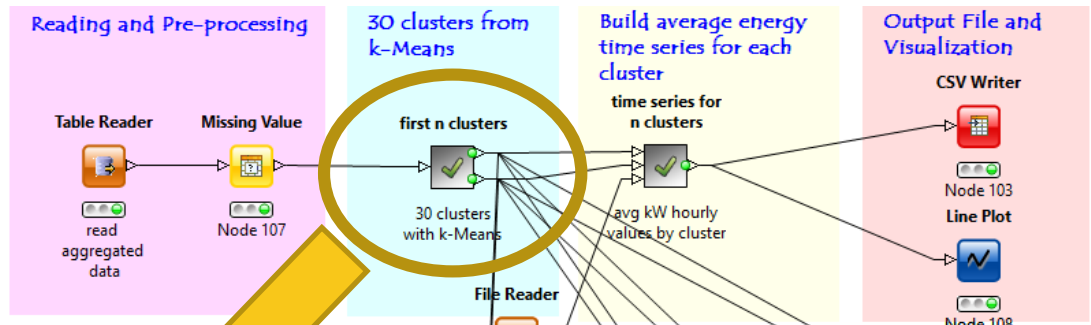
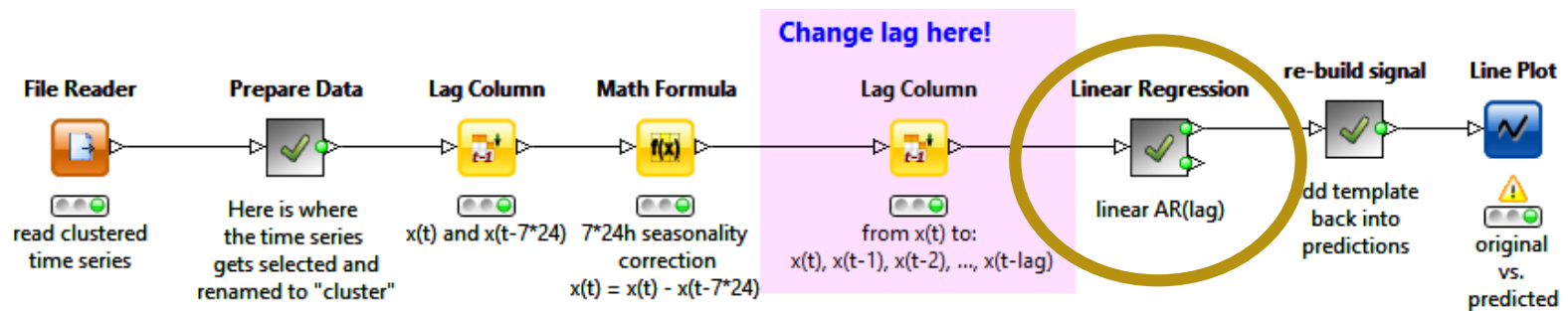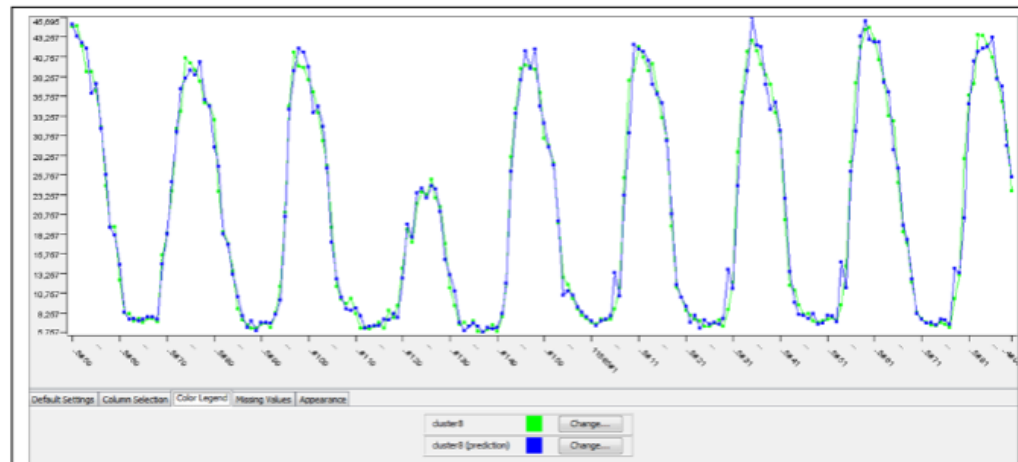- Upload a KNIME data table to Hive/Impala

# Next Steps

# Following Workflows: k-Means

K = 30

# Following Workflows: AR Model



**Change lag here!**

File Reader — Prepare Data — Lag Column — Math Formula — Lag Column — **Linear Regression** — re-build signal — Line Plot

read clustered time series

Here is where the time series gets selected and renamed to "cluster"

x(t) and x(t-7*24)

7*24h seasonality correction
x(t) = x(t) - x(t-7*24)

from x(t) to:
x(t), x(t-1), x(t-2), ..., x(t-lag)

linear AR(lag)

add template back into predictions

original vs. predicted

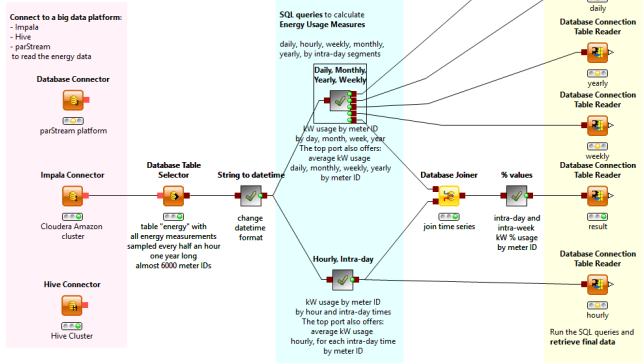*Auto-regressive model using the previous 24h\*7 as seasonality template*

- **24-hour seasonality template**: the first week of the time series is used as a template for seasonality correction
- **auto** means usage of past of the same time series for prediction. No other external time series/data used.
- **Regressive** refers to the mode used: either a linear or a polynomial regression
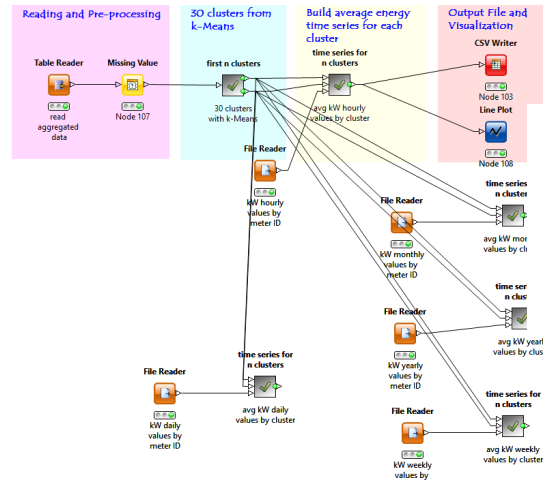
Figure 30: Original time series and predicted time series after being adjusted for weekly seasonality in green and blue respectively.

# Model Factory: Concept



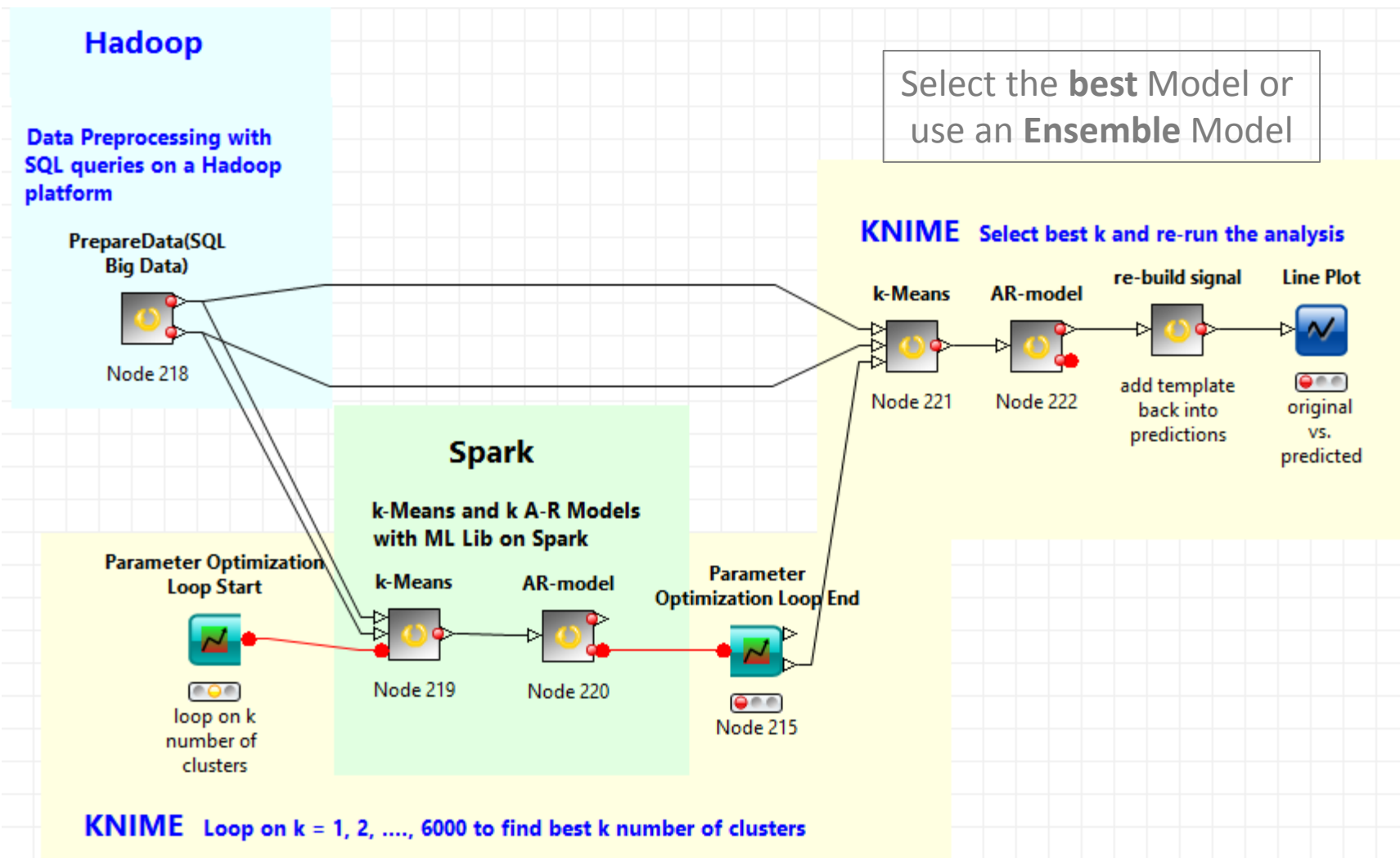Data Preparation on **Hadoop**

K-Means

A-R Model

Choosing best k for k-Means minimizing the RMS prediction error from A-R Model on **Spark**

Loop, control, and final results remain in **KNIME**

# Model Factory: Workflow

# References

- Whitepaper "KNIME opens the Doors to Big Data"
  http://www.knime.org/files/big_data_in_knime_1.pdf

- Blog Post "Integrating Big data is as Easy as 1,2,3, … 4"
  http://www.knime.org/blog/integrating-big-data-is-as-easy-as-1-2-3-4

- The Big Data Extension Product Description
  http://www.knime.org/knime-big-data-extension

Open for Innovation
KNIME

# Resources

- **KNIME** ([www.knime.org](www.knime.org))
  - **BLOG** for news, tips and tricks([www.knime.org/blog](www.knime.org/blog))
  - **FORUM** for questions and answers ([tech.knime.org/forum](tech.knime.org/forum))
  - **EXAMPLE SERVER** for example workflows
  - **LEARNING HUB** ([www.knime.org/learning-hub](www.knime.org/learning-hub))

- **KNIME TV** channel on 

- **KNIME** on  **@KNIME**

- **KNIME** on 
  **https://www.facebook.com/KNIMEanalytics**

35

# Events and Trainings

- **KNIME** (https://www.knime.org/about/events)

  - **User Training** **15-16 June** Zurich (https://www.knime.org/knime-user-training-june-2015)

  - **Developer Training** **17-18 June** Zurich (https://www.knime.org/knime-developer-training-june-2015)

- **Meetup.com**

  - **16 June**, Zurich,  Meetup on **Big Data** (http://www.meetup.com/Zurich-KNIME-Users/events/221570664/)

# Thank you!