# Application Of AI (**Knowledge Graph Embeddings**) Industrial Use Cases Formal Knowledge Integration in Machine Learning Model For Industry 4.0

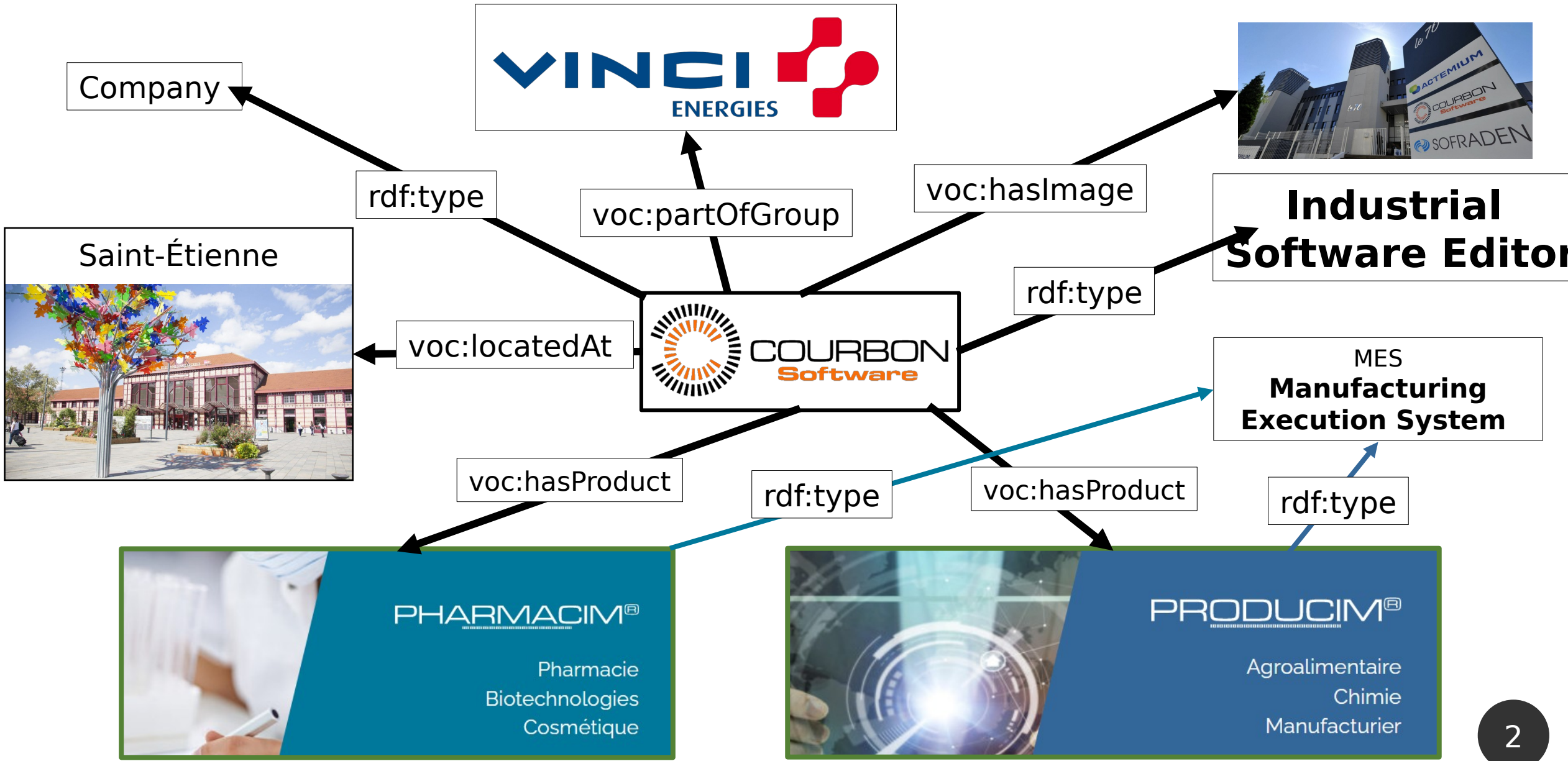**Industrial PhD Student**, Mouloud IFERROUDJENE
**Supervised by** Antoine ZIMMERMANN, Victor CHARPENAY, Thierry LAVEILLE

LIMOS

LABORATOIRE D'INFORMATIQUE,
DE MODÉLISATION ET D'OPTIMISATION DES SYSTÈMES

MINES
Saint-Étienne

anrt
ASSOCIATION NATIONALE
RECHERCHE TECHNOLOGIE

COURBON
Software

\°/ All Icons used are courtesy of flaticon.com

# COURBON Software (CSO)

MES Stands for
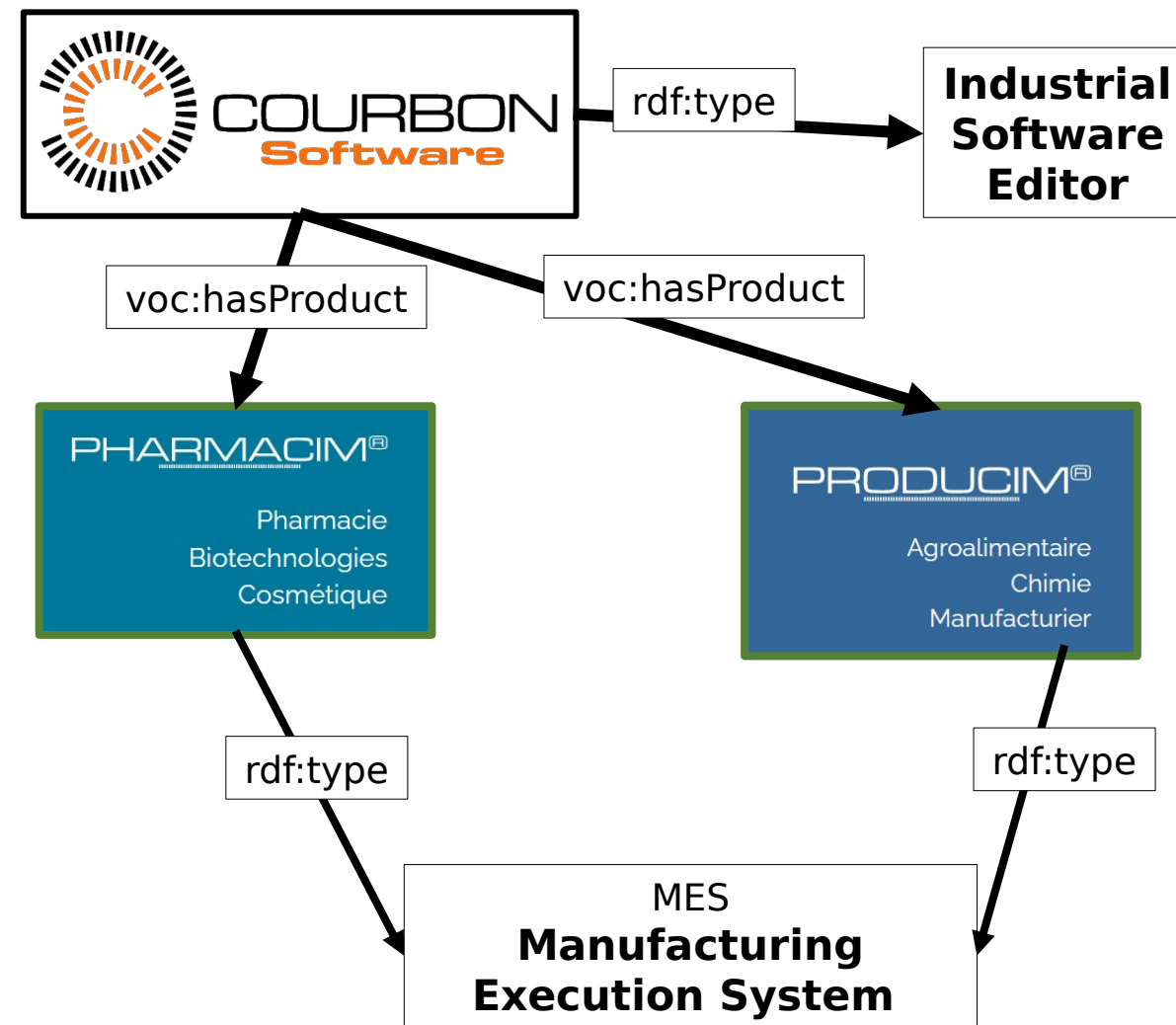➢ **Manufacturing execution system**

**MES objectives :**
➢ Ensure the proper execution of manufacturing operations
➢ Improve production efficiency

**MES Functionalities :**
➢ Product traceability
➢ Quality control
➢ Production monitoring
➢ Scheduling, Etc.

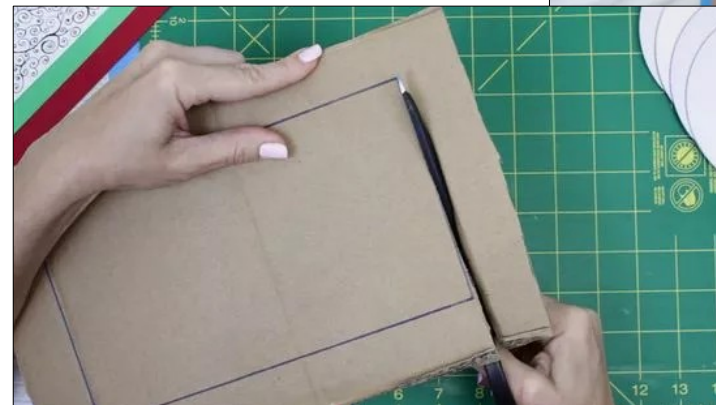**Data acquisition => Lot of DATA from different sources**

COURBON Software —rdf:type→ **Industrial Software Editor**

voc:hasProduct → PHARMACIM® Pharmacie Biotechnologies Cosmétique

voc:hasProduct → PRODUCIM® Agroalimentaire Chimie Manufacturier

PHARMACIM —rdf:type→ MES **Manufacturing Execution System**

PRODUCIM —rdf:type→ MES **Manufacturing Execution System**

# CONTEXT

# Illustration Example ( toy car production )
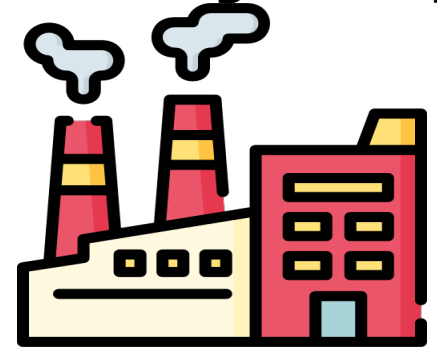
Imagine you have a toy factory where you make toy cars !!

the past, you had to do everything by hand:
Make the cars
Paint them
And package them

# Illustration Example ( car toy production )
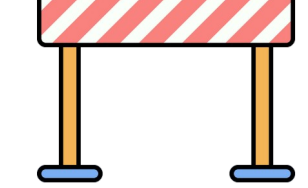
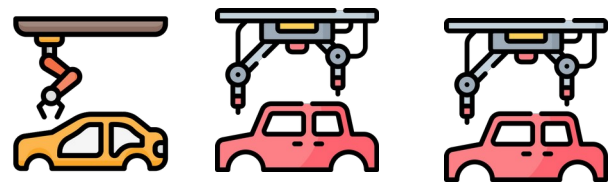**g factories (e.g., cars production)**

**The factory stop production**

**Repairing**

**Lose Money**

**ssembly line and mass production**

*Industry 2.0*

**Have a lot of machines and equipment**

**Sometimes !! machine break down**

Oops !
Oops !
Oops !
Oops !

**Machine 01**

**Machine 02**

What if we install chip to monitor our machines' status ?

*Industry 3.0*

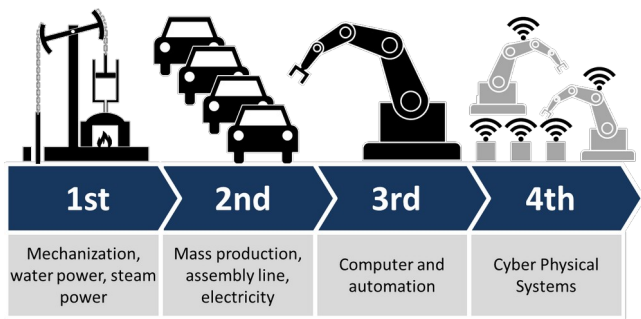**Machine 01**

**Communication**

*Industry 4.0*

**Big Data**

**Machine**

# The Fourth Industrial Revolution (Industry 4.0)



| 1st | 2nd | 3rd | 4th |
|---|---|---|---|
| Mechanization, water power, steam power | Mass production, assembly line, electricity | Computer and automation | Cyber Physical Systems |

**COURBON Software**

| Manufacturing Execution Systems (MES) | Enterprise Resource Planning (ERP) |
|---|---|

**AI applications in the industry :**

- early detection of rejects
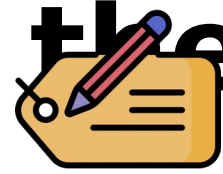- **Predictive Maintenance**
- quality control
- industrial prognosis
- etc.

**1st** — Mechanization, water power, steam power
**2nd** — Mass production, assembly line, electricity
**3rd** — Computer and automation
**4th** — Cyber Physical Systems

**Industry 4.0** + **Big data** → **Intelligent industrial system**

# Labeling of the cause of the error

- Workers labels by hand the cause of failures of every equipment, e.g.,

Noticed **High Speed** of **the conveyor**

Noticed **Low energy** of the conveyor

- Every data about the status of machine is collected
- Every issues noticed are reported and data are labled

⇒Extract insight from data to **solve** the **prevent future** the problems.

⇒The manager (or workers) gain domain-specific knowledge and expertise.

**Create Database**

| Machine N° | DateTime | Product lot N° | Issue label |
|---|---|---|---|
| M1 | 24022023T19:00:02 | 152 (Material) | 0 |
| M2 (Conveyor) | 24022023T16:02 | 12 (Toy) | Low Energy |

# Use Case – Quality Control in Cake Factory

## Example of Industrial Use Case



| Leavening Culture [ Workshop #1 ] | | Mixing, kneading [ Workshop #2 ] | | Leavening Culture [ Workshop #3 ] | | Quality Control [ Workshop #4 ] | | Cake Packaging [ Workshop #5 ] |

Description of an example **of Sequential Production-Line**

## A reminder of OiCake production line schema



**We are interested in classification problem of weight control using data collected in mixing phase (Bloc 2)**
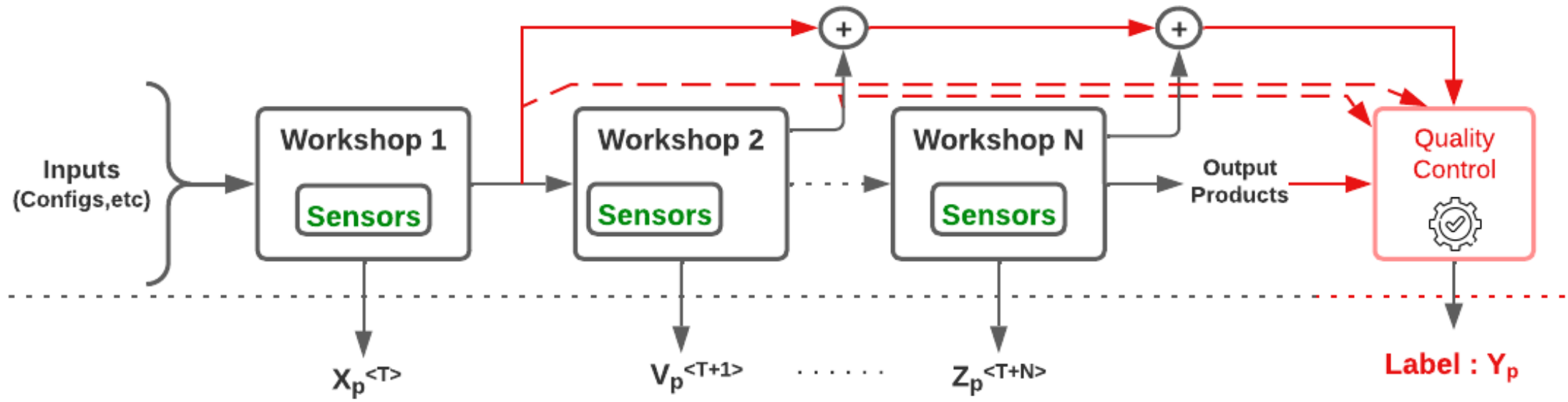
# Use Case – Quality Control in Cake Factory

## Example of Industrial Use Case
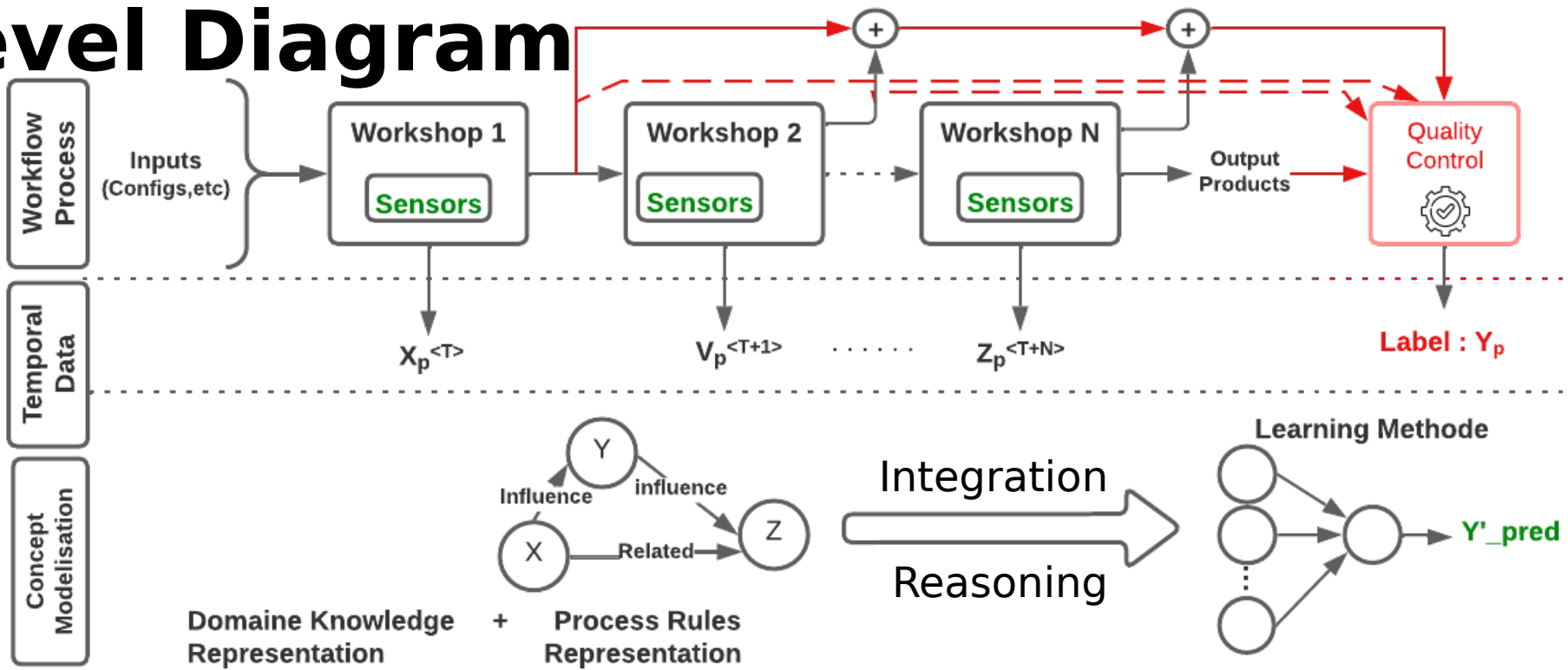
**Quality Control / Predictive Maintenance**

# Use Case – Quality Control in Cake Factory

## Workflow in Production line – Multi-Level Diagram



**Time series data**

| DATETIME | M221_TS_01 | M221_TP_01 |
|---|---|---|
| 2021-12-01 05:30:30 | 110.000000 | 6109.781280 |
| 2021-12-01 05:31:00 | 108.826667 | 6141.733192 |
| 2021-12-01 05:31:30 | 110.000000 | 6184.072883 |
| 2021-12-01 05:32:00 | 111.173333 | 6235.904952 |
| 2021-12-01 05:32:30 | 110.000000 | 6292.483363 |

**Batch information**

| | ID_LOT_PATE | ETAPE | TIME_START | TIME_STOP |
|---|---|---|---|---|
| 0 | 20211201-01 | Mélange | 2021-12-01 05:30:00 | 2021-12-01 05:44:10 |
| 1 | 20211201-01 | Pétrissage | 2021-12-01 05:44:10 | 2021-12-01 05:51:30 |
| 2 | 20211201-01 | Repos | 2021-12-01 05:59:00 | 2021-12-01 06:55:03 |
| 3 | 20211201-02 | Mélange | 2021-12-01 07:06:10 | 2021-12-01 07:20:20 |
| 4 | 20211201-02 | Pétrissage | 2021-12-01 07:20:20 | 2021-12-01 07:27:40 |

# Train ML Model for Quality Control – QICake

## Time Series data

| DATETIME | M221_TS_01 | M221_TP_01 |
|---|---|---|
| 2021-12-01 05:30:30 | 110.000000 | 6109.781280 |
| 2021-12-01 05:31:00 | 108.826667 | 6141.733192 |
| 2021-12-01 05:31:30 | 110.000000 | 6184.072883 |
| 2021-12-01 05:32:00 | 111.173333 | 6235.904952 |
| 2021-12-01 05:32:30 | 110.000000 | 6292.483363 |

## Batch information

| | ID_LOT_PATE | ETAPE | TIME_START | TIME_STOP |
|---|---|---|---|---|
| 0 | 20211201-01 | Mélange | 2021-12-01 05:30:00 | 2021-12-01 05:44:10 |
| 1 | 20211201-01 | Pétrissage | 2021-12-01 05:44:10 | 2021-12-01 05:51:30 |
| 2 | 20211201-01 | Repos | 2021-12-01 05:58:00 | 2021-12-01 06:55:00 |
| 3 | 20211201-02 | Mélange | 2021-12-01 07:06:10 | 2021-12-01 07:20:20 |
| 4 | 20211201-02 | Pétrissage | 2021-12-01 07:20:20 | 2021-12-01 07:27:40 |

## Merge Data : Batch Identifier with the temporal data (INPUT DATA)

| DATETIME | M221_TS_01 | M221_TP_01 | ID_LOT_PATE | ETAPE |
|---|---|---|---|---|
| 2021-12-01 05:30:30 | 110.000000 | 6109.781280 | 20211201-01 | Mélange |
| 2021-12-01 05:31:00 | 108.826667 | 6141.733192 | 20211201-01 | Mélange |
| 2021-12-01 05:31:30 | 110.000000 | 6184.072883 | 20211201-01 | Mélange |
| 2021-12-01 05:32:00 | 111.173333 | 6235.904952 | 20211201-01 | Mélange |
| 2021-12-01 05:32:30 | 110.000000 | 6292.483363 | 20211201-01 | Mélange |
| ... | ... | ... | ... | ... |
| 2022-03-01 12:13:30 | 118.854275 | 8742.818066 | 20220301-05 | Pétrissage |
| 2022-03-01 12:14:00 | 119.876393 | 8765.277318 | 20220301-05 | Pétrissage |
| 2022-03-01 12:14:30 | 120.888755 | 8787.267457 | 20220301-05 | Pétrissage |

## Wight Quality Control per ID_PATE (TARGET DATA)

| ID_LOT_PATE | |
|---|---|
| 20211201-01 | 1 |
| 20211201-01 | 1 |
| 20211201-01 | 0 |
| 20211201-01 | 1 |
| 20211201-02 | 0 |
| • | |
| • Only Zeros | |
| • | |
| 20211201-02 | 0 |
| 20211201-02 | 0 |

| ID_LOT_PATE | QC_wight_err |
|---|---|
| 20211201-01 | 1 |
| 20211201-02 | 0 |
| 20211201-03 | 0 |
| 20211201-04 | 1 |
| 20211201-05 | 0 |
| 20211202-01 | 1 |
| 20211202-02 | 0 |
| 20211202-03 | 1 |
| 20211202-04 | 0 |
| 20211202-05 | 0 |

# Train ML Model for Quality Control – OICake

## Model Training & Evaluation

**(1) Gradient Boosting Classifier (~ with 50 estimators)**

Train Precision : **97.2%** Train Recall: **96%**

Precision : **96.9%** Test Recall : **83.1 %**

**(2) Deep Neural Network**

Train Precision : **83.3%** Train Recall: **40.8 %**

Test Precision : **82.9%** Test Recall : **43%**

**(3) LSTM**

Train Precision: **67.1%** Train Recall: **49.5%**

Test Precision : **76.9%** Test

## Example of Industrial Use Case

**Quality Control / Predictive Maintenance**

# Workflow in Production line – Multi-Level Diagram

How to **Integrate** Domain-specific (industry 4.0) **knowledge** into **Machine Learning** to enhance its performance in downstream tasks ?

➢ **RQ1 : How** to **integrate numerical data** from **heterogeneous** observations to apply machine learning to downstream tasks? **=> C1 [ DATA INTEGRATION & HETEROGENEITY ]**

➢ **RQ2 :** How to learn **the implicit knowledge** embedded in the industrial process and reason on it using machine learning (ML) models? **=> C1 [ DATA INTEGRATION & HETEROGENEITY ] & C2 [ Explicability ]**

# Problems & Challenges

➢ **C1 [ Data integration and heterogeneity ] :** The overabundance and heterogeneity of available data **limits application of AI techniques** in the industry.

➢ **C2 [ Explicability ] :** ML/AI trained on raw data produces **black-box models** which lack **explicability**

# Proposed Approach (Overview )

## (1) Input Data



Process Data

Ontology

## (2) Integrate temporal data in knowledge graph



Workshop#w1

Workshop#w3

Workshop#QualityCntrl

Pos#Before

Pos#Before

Pos#Before

Voc:#isPartOF

R#Influence

R#Influence

R#Influence

Sen#Humidity

Sen#Temp

#Observation#01

#NumValue

Knowledge Graph

## (4) Downstream Task



Classificatio

## (3) Generate Knowledge Graph Embedding

# Industrial context - Modeling example

**ed Hierarchical Ontology Proposed for the "OICake" (cake factory) manufacturing pro**

**Membres :**
- Valentin GUIEN
- Mouloud IFERROUDJENE
- Frederic HAYEK
- Aurelien MOMBELLI

# The End

## Thank You For Your Attention



## Any Questions ?